

# 问题求解：算法与数据结构

## ( Python 版 )

1. 引言.....	4
1.1. 目标.....	4
1.2. 引子.....	4
1.3. 计算机科学是什么.....	5
1.4. 什么是程序设计.....	7
1.5. 为何要学习数据结构和抽象数据类型.....	8
1.6. 为何要学习算法.....	9
1.7. Python 入门.....	9
1.7.1. 从数据开始.....	10
1.7.2. 输入与输出.....	22
1.7.3. 控制结构.....	26
1.7.4. 异常处理.....	30
1.7.5. 定义函数.....	32
1.7.6. Python 面向对象编程：定义类.....	34
1.8. 小结.....	54
1.9. 关键词.....	54
1.10. 问题讨论.....	55
1.11. 编程练习.....	55
2. 算法分析.....	56
2.1. 目标.....	57
2.2. 什么是算法分析.....	57
2.2.1. “大 O” 表示法.....	61
2.2.2. 例子：“变位词” 判断.....	64
2.3. Python 数据结构的性能.....	69
2.3.1. 列表 List.....	69
2.3.2. 字典 Dictionary.....	74
2.4. 小结.....	76
2.5. 关键词.....	77
2.6. 问题讨论.....	77
2.7. 编程练习.....	78
3. 基本数据结构.....	78
3.1. 目标.....	79
3.2. 什么是线性结构 Linear Structure.....	79
3.3. 栈 Stack.....	79
3.3.1. 什么是栈 stack.....	79
3.3.2. 抽象数据类型 Stack.....	81
3.3.3. 用 Python 实现 Stack.....	82
3.3.4. 简单括号匹配.....	85
3.3.5. 匹配符号（通用情况）.....	87

3.3.6. 十进制数转换为二进制.....	89
3.3.7. 中缀、前缀和后缀表达式.....	92
3.4. 队列 queue.....	101
3.4.1. 什么是队列.....	101
3.4.2. 抽象数据类型 Queue.....	102
3.4.3. 在 Python 中实现 Queue.....	102
3.4.4. 模拟算法：热土豆.....	104
3.4.5. 模拟算法：打印任务.....	106
3.4.6. 主要模拟步骤.....	108
3.4.7. Python 实现.....	108
3.4.8. 讨论.....	114
3.5. 双端队列 deque.....	115
3.5.1. 什么是双端队列 deque.....	115
3.5.2. 抽象数据类型 Deque.....	115
3.5.3. 在 Python 中实现 Deque.....	116
3.5.4. “回文词”判定.....	118
3.6. 列表 List.....	120
3.6.1. 抽象数据类型无序列表 Unordered List.....	120
3.6.2. 采用链表实现无序列表.....	121
3.6.3. 抽象数据类型：有序列表 Ordered List.....	131
3.6.4. 实现有序列表.....	131
3.6.5. 链表实现算法分析.....	134
3.7. 小结.....	135
3.8. 关键词.....	135
3.9. 问题讨论.....	136
3.10. 编程练习.....	136
4. 递归 Recursion.....	139
4.1. 目标.....	140
4.2. 什么是递归.....	140
4.2.1. 计算数列表的和.....	140
4.2.2. 递归的三法则.....	142
4.2.3. 将整数转换为字符串形式的任意进制表示.....	143
4.3. 栈帧：实现递归.....	146
4.4. 图示递归.....	148
4.4.1. 谢尔宾斯基三角形 Sierpinski Triangle.....	152
4.5. 复杂递归问题.....	155
4.5.1. 河内塔问题 Towers of Hanoi.....	155
4.6. 探索迷宫.....	158
4.7. 动态规划.....	166
4.8. 小结.....	172
4.9. 关键词.....	172
4.10. 问题讨论.....	173
4.11. 词汇表.....	173
4.12. 编程练习.....	174

5. 排序与搜索.....	176
5.1. 目标.....	176
5.2. 搜索.....	176
5.2.1. 顺序搜索.....	176
5.2.2. 二分法搜索.....	179
5.2.3. 散列.....	183
5.3. 排序.....	194
5.3.1. 冒泡排序.....	194
5.3.2. 选择排序.....	197
5.3.3. 插入排序.....	199
5.3.4. 希尔排序.....	201
5.3.5. 归并排序.....	204
5.3.6. 快速排序.....	207
5.4. 小结.....	212
5.5. 关键词.....	212
5.6. 问题讨论.....	213
5.7. 编程练习.....	214
6. 树和树算法.....	215
6.1. 目标.....	216
6.2. 树的例子.....	216
6.3. 术语表与定义.....	217
6.3.1. 术语表.....	217
6.3.2. 定义.....	217
6.4. 实现.....	217
6.4.1. “列表的列表”表示树.....	217
6.4.2. 节点和引用.....	217
6.5. 二叉堆 Binary Heap 实现的优先队列.....	217
6.5.1. 二叉堆操作.....	217
6.5.2. 二叉堆实现.....	217
6.6. 二叉树应用.....	217
6.7. 树遍历.....	217
6.8. 二叉搜索树.....	217
6.8.1. 搜索树操作.....	217
6.8.2. 搜索树实现.....	217
6.8.3. 搜索树分析.....	217
6.9. 小结.....	217
6.10. 关键词.....	218
6.11. 问题讨论.....	218
6.12. 编程练习.....	218
7. 图和图算法.....	218
7.1. 目标.....	218
7.2. 词汇表及定义.....	218
7.3. 图抽象数据类型.....	218
7.4. 邻接矩阵.....	218

7.5. 邻接表.....	218
7.6. 实现.....	218
7.7. Word Ladder 词梯问题.....	218
7.7.1. 建立 Word Ladder 图.....	218
7.7.2. 实现广度优先搜索.....	218
7.7.3. 广度优先搜索分析.....	218
7.8. 骑士周游问题.....	218
7.8.1. 建立骑士周游图.....	219
7.8.2. 实现骑士周游.....	219
7.8.3. 骑士周游分析.....	219
7.8.4. 通用深度优先搜索.....	219
7.8.5. 深度优先分析.....	219
7.9. 拓扑排序.....	219
7.10. 强连通分支.....	219
7.11. 最短路径问题.....	219
7.11.1. Dijkstra 算法.....	219
7.11.2. Dijkstra 算法分析.....	219
7.12. Prim 最小生成树算法.....	219
7.13. 小结.....	219
7.14. 关键词.....	219
7.15. 问题讨论.....	219
7.16. 编程练习.....	219

# 1. 引言

## 1.1. 目标

- 关于计算机科学、程序设计和问题求解的基本概念；
- 什么是“抽象”，及抽象在问题求解过程中的作用；
- 什么是“抽象数据类型”，及其实现；
- Python 程序设计语言入门

## 1.2. 引子

The way we think about programming has undergone many changes in the years since the first electronic computers required patch cables and switches to convey instructions from human to machine. As is the case with many aspects of society, changes in computing technology provide computer scientists with a growing number of tools and platforms on which to practice their craft. Advances such as faster processors, high-speed networks, and large memory capacities have created a spiral of complexity through which computer scientists must navigate. Throughout all of

this rapid evolution, a number of basic principles have remained constant. The science of computing is concerned with using computers to solve problems.

You have no doubt spent considerable time learning the basics of problem-solving and hopefully feel confident in your ability to take a problem statement and develop a solution. You have also learned that writing computer programs is often hard. The complexity of large problems and the corresponding complexity of the solutions can tend to overshadow the fundamental ideas related to the problem-solving process.

This chapter emphasizes two important areas for the rest of the text. First, it reviews the framework within which computer science and the study of algorithms and data structures must fit, in particular, the reasons why we need to study these topics and how understanding these topics helps us to become better problem solvers. Second, we review the Python programming language. Although we cannot provide a detailed, exhaustive reference, we will give examples and explanations for the basic constructs and ideas that will occur throughout the remaining chapters.

### 1.3. 计算机科学是什么

Computer science is often difficult to define. This is probably due to the unfortunate use of the word “computer” in the name. As you are perhaps aware, computer science is not simply the study of computers. Although computers play an important supporting role as a tool in the discipline, they are just that – tools.

Computer science is the study of problems, problem-solving, and the solutions that come out of the problem-solving process. Given a problem, a computer scientist’s goal is to develop an algorithm, a step-by-step list of instructions for solving any instance of the problem that might arise. Algorithms are finite processes that if followed will solve the problem. Algorithms are solutions.

Computer science can be thought of as the study of algorithms. However, we must be careful to include the fact that some problems may not have a solution. Although proving this statement is beyond the scope of this text, the fact that some problems cannot be solved is important for those who study computer science. We can fully define computer science, then, by including both types of problems and stating that computer science is the study of solutions to problems as well as the study of problems with no solutions.

It is also very common to include the word computable when describing problems and solutions. We say that a problem is computable if an algorithm exists for solving it. An alternative definition for computer science, then, is to say that computer science is the study of problems that are and that are not computable, the study of the existence and the nonexistence of algorithms. In any case, you will note that the word “computer” did not come up at all. Solutions are considered independent from the machine.

Computer science, as it pertains to the problem-solving process itself, is also the study of

abstraction. Abstraction allows us to view the problem and solution in such a way as to separate the so-called logical and physical perspectives. The basic idea is familiar to us in a common example.

Consider the automobile that you may have driven to school or work today. As a driver, a user of the car, you have certain interactions that take place in order to utilize the car for its intended purpose. You get in, insert the key, start the car, shift, brake, accelerate, and steer in order to drive. From an abstraction point of view, we can say that you are seeing the logical perspective of the automobile. You are using the functions provided by the car designers for the purpose of transporting you from one location to another. These functions are sometimes also referred to as the interface.

On the other hand, the mechanic who must repair your automobile takes a very different point of view. She not only knows how to drive but must know all of the details necessary to carry out all the functions that we take for granted. She needs to understand how the engine works, how the transmission shifts gears, how temperature is controlled, and so on. This is known as the physical perspective, the details that take place “under the hood.”

The same thing happens when we use computers. Most people use computers to write documents, send and receive email, surf the web, play music, store images, and play games without any knowledge of the details that take place to allow those types of applications to work. They view computers from a logical or user perspective. Computer scientists, programmers, technology support staff, and system administrators take a very different view of the computer. They must know the details of how operating systems work, how network protocols are configured, and how to code various scripts that control function. They must be able to control the low-level details that a user simply assumes.

The common point for both of these examples is that the user of the abstraction, sometimes also called the client, does not need to know the details as long as the user is aware of the way the interface works. This interface is the way we as users communicate with the underlying complexities of the implementation. As another example of abstraction, consider the Python math module. Once we import the module, we can perform computations such as

```
>>> import math
>>> math.sqrt(16)
4.0
>>>
```

This is an example of procedural abstraction. We do not necessarily know how the square root is being calculated, but we know what the function is called and how to use it. If we perform the import correctly, we can assume that the function will provide us with the correct results. We know that someone implemented a solution to the square root problem but we only need to know how to use it. This is sometimes referred to as a “black box” view of a process. We simply describe the interface: the name of the function, what is needed (the parameters), and what will be

returned. The details are hidden inside (see Figure 1).

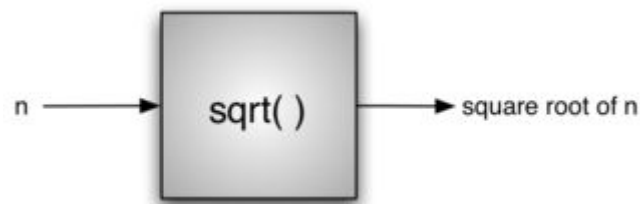


图 1 Procedural Abstraction

## 1.4. 什么是程序设计

Programming is the process of taking an algorithm and encoding it into a notation, a programming language, so that it can be executed by a computer. Although many programming languages and many different types of computers exist, the important first step is the need to have the solution. Without an algorithm there can be no program.

Computer science is not the study of programming. Programming, however, is an important part of what a computer scientist does. Programming is often the way that we create a representation for our solutions. Therefore, this language representation and the process of creating it becomes a fundamental part of the discipline.

Algorithms describe the solution to a problem in terms of the data needed to represent the problem instance and the set of steps necessary to produce the intended result. Programming languages must provide a notational way to represent both the process and the data. To this end, languages provide control constructs and data types.

Control constructs allow algorithmic steps to be represented in a convenient yet unambiguous way. At a minimum, algorithms require constructs that perform sequential processing, selection for decision-making, and iteration for repetitive control. As long as the language provides these basic statements, it can be used for algorithm representation.

All data items in the computer are represented as strings of binary digits. In order to give these strings meaning, we need to have data types. Data types provide an interpretation for this binary data so that we can think about the data in terms that make sense with respect to the problem being solved. These low-level, built-in data types (sometimes called the primitive data types) provide the building blocks for algorithm development.

For example, most programming languages provide a data type for integers. Strings of binary digits in the computer's memory can be interpreted as integers and given the typical meanings that we commonly associate with integers (e.g. 23, 654, and -19). In addition, a data type also provides a description of the operations that the data items can participate in. With integers, operations such as addition, subtraction, and multiplication are common. We have come to expect that numeric types of data can participate in these arithmetic operations.

The difficulty that often arises for us is the fact that problems and their solutions are very complex. These simple, language-provided constructs and data types, although certainly sufficient to represent complex solutions, are typically at a disadvantage as we work through the problem-solving process. We need ways to control this complexity and assist with the creation of solutions.

## 1.5. 为何要学习数据结构和抽象数据类型

To manage the complexity of problems and the problem-solving process, computer scientists use abstractions to allow them to focus on the “big picture” without getting lost in the details. By creating models of the problem domain, we are able to utilize a better and more efficient problem-solving process. These models allow us to describe the data that our algorithms will manipulate in a much more consistent way with respect to the problem itself.

Earlier, we referred to procedural abstraction as a process that hides the details of a particular function to allow the user or client to view it at a very high level. We now turn our attention to a similar idea, that of data abstraction. An abstract data type, sometimes abbreviated ADT, is a logical description of how we view the data and the operations that are allowed without regard to how they will be implemented. This means that we are concerned only with what the data is representing and not with how it will eventually be constructed. By providing this level of abstraction, we are creating an encapsulation around the data. The idea is that by encapsulating the details of the implementation, we are hiding them from the user’s view. This is called **information hiding**.

Figure 2 shows a picture of what an abstract data type is and how it operates. The user interacts with the interface, using the operations that have been specified by the abstract data type. The abstract data type is the shell that the user interacts with. The implementation is hidden one level deeper. The user is not concerned with the details of the implementation.

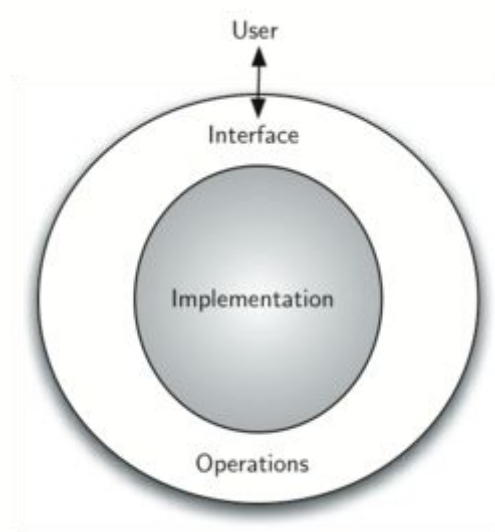


图 2 Abstract Data Type



The implementation of an abstract data type, often referred to as a data structure, will require that we provide a physical view of the data using some collection of programming constructs and primitive data types. As we discussed earlier, the separation of these two perspectives will allow us to define the complex data models for our problems without giving any indication as to the details of how the model will actually be built. This provides an implementation-independent view of the data. Since there will usually be many different ways to implement an abstract data type, this implementation independence allows the programmer to switch the details of the implementation without changing the way the user of the data interacts with it. The user can remain focused on the problem-solving process.

## 1.6. 为何要学习算法

Computer scientists learn by experience. We learn by seeing others solve problems and by solving problems by ourselves. Being exposed to different problem-solving techniques and seeing how different algorithms are designed helps us to take on the next challenging problem that we are given. By considering a number of different algorithms, we can begin to develop pattern recognition so that the next time a similar problem arises, we are better able to solve it.

Algorithms are often quite different from one another. Consider the example of `sqrt` seen earlier. It is entirely possible that there are many different ways to implement the details to compute the square root function. One algorithm may use many fewer resources than another. One algorithm might take 10 times as long to return the result as the other. We would like to have some way to compare these two solutions. Even though they both work, one is perhaps “better” than the other. We might suggest that one is more efficient or that one simply works faster or uses less memory. As we study algorithms, we can learn analysis techniques that allow us to compare and contrast solutions based solely on their own characteristics, not the characteristics of the program or computer used to implement them.

In the worst case scenario, we may have a problem that is intractable, meaning that there is no algorithm that can solve the problem in a realistic amount of time. It is important to be able to distinguish between those problems that have solutions, those that do not, and those where solutions exist but require too much time or other resources to work reasonably.

There will often be trade-offs that we will need to identify and decide upon. As computer scientists, in addition to our ability to solve problems, we will also need to know and understand solution evaluation techniques. In the end, there are often many ways to solve a problem. Finding a solution and then deciding whether it is a good one are tasks that we will do over and over again.

## 1.7. Python 入门

In this section, we will review the programming language Python and also provide some more detailed examples of the ideas from the previous section. If you are new to Python or find that you

need more information about any of the topics presented, we recommend that you consult a resource such as the **Python Language Reference** or a **Python Tutorial**. Our goal here is to reacquaint you with the language and also reinforce some of the concepts that will be central to later chapters.

Python is a modern, easy-to-learn, object-oriented programming language. It has a powerful set of built-in data types and easy-to-use control constructs. Since Python is an interpreted language, it is most easily reviewed by simply looking at and describing interactive sessions. You should recall that the interpreter displays the familiar `>>>` prompt and then evaluates the Python construct that you provide. For example,

```
>>> print("Algorithms and Data Structures")
Algorithms and Data Structures
>>>
```

shows the prompt, the `print` function, the result, and the next prompt.

### 1.7.1. 从数据开始

We stated above that Python supports the object-oriented programming paradigm. This means that Python considers data to be the focal point of the problem-solving process. In Python, as well as in any other object-oriented programming language, we define a class to be a description of what the data look like (the state) and what the data can do (the behavior). Classes are analogous to abstract data types because a user of a class only sees the state and behavior of a data item. Data items are called objects in the object-oriented paradigm. An object is an instance of a class.

#### 1.7.1.1. 内建原子数据类型

We will begin our review by considering the atomic data types. Python has two main built-in numeric classes that implement the integer and floating point data types. These Python classes are called `int` and `float`. The standard arithmetic operations, `+`, `-`, `*`, `/`, and `**` (exponentiation), can be used with parentheses forcing the order of operations away from normal operator precedence. Other very useful operations are the remainder (modulo) operator, `%`, and integer division, `//`. Note that when two integers are divided, the result is a floating point. The integer division operator returns the integer portion of the quotient by truncating any fractional part.

```
print 2+3*4
print (2+3)*4
print 2**10
print 6/3
print 7/3
print 7//3
```

```
print 7%3
print 3/6
print 3//6
print 3%6
print 2**100
```

The boolean data type, implemented as the Python `bool` class, will be quite useful for representing truth values. The possible state values for a boolean object are `True` and `False` with the standard boolean operators, `and`, `or`, and `not`.

```
>>> True
True
>>> False
False
>>> False or True
True
>>> not (False or True)
False
>>> True and True
True
```

代码 1 Basic Arithmetic Operators (intro\_1)

Boolean data objects are also used as results for comparison operators such as equality (`==`) and greater than (`>`). In addition, relational operators and logical operators can be combined together to form complex logical questions. Table 1 shows the relational and logical operators with examples shown in the session that follows.

Operation Name	Operator	Explanation
less than	<code>&lt;</code>	Less than operator
greater than	<code>&gt;</code>	Greater than operator
less than or equal	<code>&lt;=</code>	Less than or equal to operator
greater than or equal	<code>&gt;=</code>	Greater than or equal to operator
equal	<code>==</code>	Equality operator
not equal	<code>!=</code>	Not equal operator
logical and	<code>and</code>	Both operands True for result to be True
logical or	<code>or</code>	One or the other operand is True for the result to be True
logical not	<code>not</code>	Negates the truth value, False becomes True, True becomes False

表格 1 Relational and Logical Operators

```
print(5==10)
```

```
print(10 > 5)
print((5 >= 1) and (5 <= 10))
```

代码 2 Basic Relational and Logical Operators (intro\_2)

Identifiers are used in programming languages as names. In Python, identifiers start with a letter or an underscore (`_`), are case sensitive, and can be of any length. Remember that it is always a good idea to use names that convey meaning so that your program code is easier to read and understand.

A Python variable is created when a name is used for the first time on the left-hand side of an assignment statement. Assignment statements provide a way to associate a name with a value. The variable will hold a reference to a piece of data and not the data itself. Consider the following session:

```
>>> theSum = 0
>>> theSum
0
>>> theSum = theSum + 1
>>> theSum
1
>>> theSum = True
>>> theSum
True
```

The assignment statement `theSum = 0` creates a variable called `theSum` and lets it hold the reference to the data object `0` (see Figure 3). In general, the right-hand side of the assignment statement is evaluated and a reference to the resulting data object is “assigned” to the name on the left-hand side. At this point in our example, the type of the variable is integer as that is the type of the data currently being referred to by `thesum`. If the type of the data changes (see Figure 4), as shown above with the boolean value `True`, so does the type of the variable (`thesum` is now of the type boolean). The assignment statement changes the reference being held by the variable. This is a dynamic characteristic of Python. The same variable can refer to many different types of data.



图 3 Variables Hold References to Data Objects

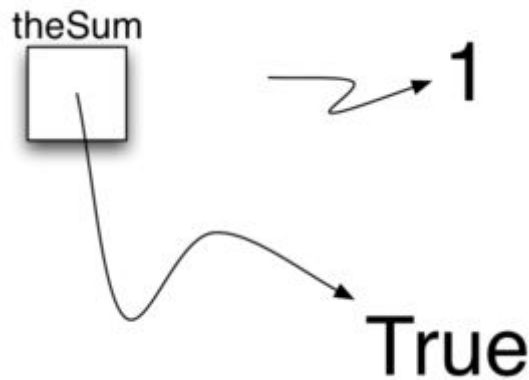


图 4 Assignment Changes the Reference

### 1.7.1.2. 内建集合数据类型

In addition to the numeric and boolean classes, Python has a number of very powerful built-in collection classes. Lists, strings, and tuples are ordered collections that are very similar in general structure but have specific differences that must be understood for them to be used properly. Sets and dictionaries are unordered collections.

A **list** is an ordered collection of zero or more references to Python data objects. Lists are written as comma-delimited values enclosed in square brackets. The empty list is simply `[]`. Lists are heterogeneous, meaning that the data objects need not all be from the same class and the collection can be assigned to a variable as below. The following fragment shows a variety of Python data objects in a list.

```
>>> [1,3,True,6.5]
[1, 3, True, 6.5]
>>> myList = [1,3,True,6.5]
>>> myList
[1, 3, True, 6.5]
```

Note that when Python evaluates a list, the list itself is returned. However, in order to remember the list for later processing, its reference needs to be assigned to a variable.

Since lists are considered to be sequentially ordered, they support a number of operations that can be applied to any Python sequence. Table 2 reviews these operations and the following session gives examples of their use.

Operation Name	Operator	Explanation
----------------	----------	-------------

indexing	[ ]	Access an element of a sequence
concatenation	+	Combine sequences together
repetition	*	Concatenate a repeated number of times
membership	in	Ask whether an item is in a sequence
length	len	Ask the number of items in the sequence
slicing	[ : ]	Extract a part of a sequence

表格 2 Operations on Any Sequence in Python

Note that the indices for lists (sequences) start counting with 0. The slice operation, `myList[1:3]`, returns a list of items starting with the item indexed by 1 up to but not including the item indexed by 3.

Sometimes, you will want to initialize a list. This can quickly be accomplished by using repetition. For example,

```
>>> myList = [0] * 6
>>> myList
[0, 0, 0, 0, 0, 0]
```

One very important aside relating to the repetition operator is that the result is a repetition of references to the data objects in the sequence. This can best be seen by considering the following session:

```
myList = [1,2,3,4]
A = [myList]*3
print(A)
myList[2]=45
print(A)
```

代码 3 Repetition of References (intro\_3)

The variable `A` holds a collection of three references to the original list called `myList`. Note that a change to one element of `myList` shows up in all three occurrences in `A`.

Lists support a number of methods that will be used to build data structures. Table 3 provides a summary. Examples of their use follow.

Method Name	Use	Explanation
<code>append</code>	<code>alist.append(item)</code>	Adds a new item to the end of a list
<code>insert</code>	<code>alist.insert(i,item)</code>	Inserts an item at the <i>i</i> th position in a list
<code>pop</code>	<code>alist.pop()</code>	Removes and returns the last item in a list
<code>pop</code>	<code>alist.pop(i)</code>	Removes and returns the <i>i</i> th item in a list
<code>sort</code>	<code>alist.sort()</code>	Modifies a list to be sorted

<code>reverse</code>	<code>alist.reverse()</code>	Modifies a list to be in reverse order
<code>del</code>	<code>del alist[i]</code>	Deletes the item in the ith position
<code>index</code>	<code>alist.index(item)</code>	Returns the index of the first occurrence of item
<code>count</code>	<code>alist.count(item)</code>	Returns the number of occurrences of item
<code>remove</code>	<code>alist.remove(item)</code>	Removes the first occurrence of item

表格 3 Methods Provided by Lists in Python

```
myList = [1024, 3, True, 6.5]
myList.append(False)
print(myList)
myList.insert(2,4.5)
print(myList)
print(myList.pop())
print(myList)
print(myList.pop(1))
print(myList)
myList.pop(2)
print(myList)
myList.sort()
print(myList)
myList.reverse()
print(myList)
print(myList.count(6.5))
print(myList.index(4.5))
myList.remove(6.5)
print(myList)
del myList[0]
print(myList)
```

代码 4 Examples of List Methods (intro\_5)

You can see that some of the methods, such as `pop`, return a value and also modify the list. Others, such as `reverse`, simply modify the list with no return value. `pop` will default to the end of the list but can also remove and return a specific item. The index range starting from 0 is again used for these methods. You should also notice the familiar “dot” notation for asking an object to invoke a method. `myList.append(False)` can be read as “ask the object `myList` to perform its `append` method and send it the value `False`.” Even simple data objects such as integers can invoke methods in this way.

```
>>> (54).__add__(21)
75
>>>
```

In this fragment we are asking the integer object `54` to execute its `add` method (called `__add__` in

Python) and passing it `21` as the value to add. The result is the sum, `75`. Of course, we usually write this as `54+21`. We will say much more about these methods later in this section.

One common Python function that is often discussed in conjunction with lists is the `range` function. `range` produces a range object that represents a sequence of values. By using the `list` function, it is possible to see the value of the range object as a list. This is illustrated below.

```
>>> range(10)
range(0, 10)
>>> list(range(10))
[0, 1, 2, 3, 4, 5, 6, 7, 8, 9]
>>> range(5,10)
range(5, 10)
>>> list(range(5,10))
[5, 6, 7, 8, 9]
>>> list(range(5,10,2))
[5, 7, 9]
>>> list(range(10,1,-1))
[10, 9, 8, 7, 6, 5, 4, 3, 2]
>>>
```

The range object represents a sequence of integers. By default, it will start with 0. If you provide more parameters, it will start and end at particular points and can even skip items. In our first example, `range(10)`, the sequence starts with 0 and goes up to but does not include 10. In our second example, `range(5,10)` starts at 5 and goes up to but not including 10. `range(5,10,2)` performs similarly but skips by twos (again, 10 is not included).

**Strings** are sequential collections of zero or more letters, numbers and other symbols. We call these letters, numbers and other symbols characters. Literal string values are differentiated from identifiers by using quotation marks (either single or double).

```
>>> "David"
'David'
>>> myName = "David"
>>> myName[3]
'i'
>>> myName*2
'DavidDavid'
>>> len(myName)
5
>>>
```

Since strings are sequences, all of the sequence operations described above work as you would



expect. In addition, strings have a number of methods, some of which are shown in Table 4. For example,

```
>>> myName
'David'
>>> myName.upper()
'DAVID'
>>> myName.center(10)
'  David  '
>>> myName.find('v')
2
>>> myName.split('v')
['Da', 'id']
```

Of these, `split` will be very useful for processing data. `split` will take a string and return a list of strings using the split character as a division point. In the example, `v` is the division point. If no division is specified, the `split` method looks for whitespace characters such as tab, newline and space.

Method Name	Use	Explanation
<code>center</code>	<code>astring.center(w)</code>	Returns a string centered in a field of size <code>w</code>
<code>count</code>	<code>astring.count(item)</code>	Returns the number of occurrences of <code>item</code> in the string
<code>ljust</code>	<code>astring.ljust(w)</code>	Returns a string left-justified in a field of size <code>w</code>
<code>lower</code>	<code>astring.lower()</code>	Returns a string in all lowercase
<code>rjust</code>	<code>astring.rjust(w)</code>	Returns a string right-justified in a field of size <code>w</code>
<code>find</code>	<code>astring.find(item)</code>	Returns the index of the first occurrence of <code>item</code>
<code>split</code>	<code>astring.split(schar)</code>	Splits a string into substrings at <code>schar</code>

表格 4 Methods Provided by Strings in Python

A major difference between lists and strings is that lists can be modified while strings cannot. This is referred to as **mutability**. Lists are mutable; strings are immutable. For example, you can change an item in a list by using indexing and assignment. With a string that change is not allowed.

```
>>> myList
[1, 3, True, 6.5]
>>> myList[0]=2**10
>>> myList
[1024, 3, True, 6.5]
>>>
>>> myName
```

```
'David'
>>> myName[0]='X'

Traceback (most recent call last):
  File "<pyshell#84>", line 1, in -toplevel-
    myName[0]='X'
TypeError: object doesn't support item assignment
>>>
```

Tuples are very similar to lists in that they are heterogeneous sequences of data. The difference is that a tuple is immutable, like a string. A tuple cannot be changed. Tuples are written as comma-delimited values enclosed in parentheses. As sequences, they can use any operation described above. For example,

```
>>> myTuple = (2,True,4.96)
>>> myTuple
(2, True, 4.96)
>>> len(myTuple)
3
>>> myTuple[0]
2
>>> myTuple * 3
(2, True, 4.96, 2, True, 4.96, 2, True, 4.96)
>>> myTuple[0:2]
(2, True)
>>>
```

However, if you try to change an item in a tuple, you will get an error. Note that the error message provides location and reason for the problem.

```
>>> myTuple[1]=False

Traceback (most recent call last):
  File "<pyshell#137>", line 1, in -toplevel-
    myTuple[1]=False
TypeError: object doesn't support item assignment
>>>
```

A set is an unordered collection of zero or more immutable Python data objects. Sets do not allow duplicates and are written as comma-delimited values enclosed in curly braces. The empty set is represented by `set()`. Sets are heterogeneous, and the collection can be assigned to a variable as below.

```
>>> {3,6,"cat",4.5,False}
{False, 4.5, 3, 6, 'cat'}
>>> mySet = {3,6,"cat",4.5,False}
>>> mySet
{False, 4.5, 3, 6, 'cat'}
>>>
```

Even though sets are not considered to be sequential, they do support a few of the familiar operations presented earlier. Table 5 reviews these operations and the following session gives examples of their use.

Operation Name	Operator	Explanation
membership	<code>in</code>	Set membership
length	<code>len</code>	Returns the cardinality of the set
<code> </code>	<code>aset   otherset</code>	Returns a new set with all elements from both sets
<code>&amp;</code>	<code>aset &amp; otherset</code>	Returns a new set with only those elements common to both sets
<code>-</code>	<code>aset - otherset</code>	Returns a new set with all items from the first set not in second
<code>&lt;=</code>	<code>aset &lt;= otherset</code>	Asks whether all elements of the first set are in the second

表格 5 Operations on a Set in Python

```
>>> mySet
{False, 4.5, 3, 6, 'cat'}
>>> len(mySet)
5
>>> False in mySet
True
>>> "dog" in mySet
False
>>>
```

Sets support a number of methods that should be familiar to those who have worked with them in a mathematics setting. Table 6 provides a summary. Examples of their use follow. Note that `union`, `intersection`, `issubset`, and `difference` all have operators that can be used as well.

Method Name	Use	Explanation
<code>union</code>	<code>aset.union(otherset)</code>	Returns a new set with all elements from

		both sets
intersection	aset.intersection(otherse t)	Returns a new set with only those elements common to both sets
difference	aset.difference(otherset)	Returns a new set with all items from first set not in second
issubset	aset.issubset(otherset)	Asks whether all elements of one set are in the other
add	aset.add(item)	Adds item to the set
remove	aset.remove(item)	Removes item from the set
pop	aset.pop()	Removes an arbitrary element from the set
clear	aset.clear()	Removes all elements from the set

表格 6 Methods Provided by Sets in Python

```

>>> mySet
{False, 4.5, 3, 6, 'cat'}
>>> yourSet = {99,3,100}
>>> mySet.union(yourSet)
{False, 4.5, 3, 100, 6, 'cat', 99}
>>> mySet | yourSet
{False, 4.5, 3, 100, 6, 'cat', 99}
>>> mySet.intersection(yourSet)
{3}
>>> mySet & yourSet
{3}
>>> mySet.difference(yourSet)
{False, 4.5, 6, 'cat'}
>>> mySet - yourSet
{False, 4.5, 6, 'cat'}
>>> {3,100}.issubset(yourSet)
True
>>> {3,100}<=yourSet
True
>>> mySet.add("house")
>>> mySet
{False, 4.5, 3, 6, 'house', 'cat'}
>>> mySet.remove(4.5)
>>> mySet
{False, 3, 6, 'house', 'cat'}
>>> mySet.pop()
False
>>> mySet
{3, 6, 'house', 'cat'}
>>> mySet.clear()

```

```
>>> mySet
set()
>>>
```

Our final Python collection is an unordered structure called a **dictionary**. Dictionaries are collections of associated pairs of items where each pair consists of a key and a value. This key-value pair is typically written as key:value. Dictionaries are written as comma-delimited key:value pairs enclosed in curly braces. For example,

```
>>> capitals = {'Iowa':'DesMoines','Wisconsin':'Madison'}
>>> capitals
{'Wisconsin': 'Madison', 'Iowa': 'DesMoines'}
>>>
```

We can manipulate a dictionary by accessing a value via its key or by adding another key-value pair. The syntax for access looks much like a sequence access except that instead of using the index of the item we use the key value. To add a new value is similar.

```
capitals = {'Iowa':'DesMoines','Wisconsin':'Madison'}
print(capitals['Iowa'])
capitals['Utah']='SaltLakeCity'
print(capitals)
capitals['California']='Sacramento'
print(len(capitals))
for k in capitals:
    print(capitals[k]," is the capital of ", k)
```

代码 5 Using a Dictionary (intro\_7)

It is important to note that the dictionary is maintained in no particular order with respect to the keys. The first pair added ('Utah': 'SaltLakeCity') was placed first in the dictionary and the second pair added ('California': 'Sacramento') was placed last. The placement of a key is dependent on the idea of “hashing,” which will be explained in more detail in Chapter 4. We also show the length function performing the same role as with previous collections.

Dictionaries have both methods and operators. Table 7 and Table 8 describe them, and the session shows them in action. The **keys**, **values**, and **items** methods all return objects that contain the values of interest. You can use the **list** function to convert them to lists. You will also see that there are two variations on the **get** method. If the key is not present in the dictionary, **get** will return **None**. However, a second, optional parameter can specify a return value instead.

Operator	Use	Explanation
[ ]	myDict[k]	Returns the value associated with k, otherwise its an error

<code>in</code>	<code>key in adict</code>	Returns <code>True</code> if key is in the dictionary, <code>False</code> otherwise
<code>del</code>	<code>del adict[key]</code>	Removes the entry from the dictionary

表格 7 Operators Provided by Dictionaries in Python

```
>>> phoneext={'david':1410,'brad':1137}
>>> phoneext
{'brad': 1137, 'david': 1410}
>>> phoneext.keys()
dict_keys(['brad', 'david'])
>>> list(phoneext.keys())
['brad', 'david']
>>> phoneext.values()
dict_values([1137, 1410])
>>> list(phoneext.values())
[1137, 1410]
>>> phoneext.items()
dict_items([('brad', 1137), ('david', 1410)])
>>> list(phoneext.items())
[('brad', 1137), ('david', 1410)]
>>> phoneext.get("kent")
>>> phoneext.get("kent","NO ENTRY")
'NO ENTRY'
>>>
```

Method Name	Use	Explanation
<code>keys</code>	<code>adict.keys()</code>	Returns the keys of the dictionary in a <code>dict_keys</code> object
<code>values</code>	<code>adict.values()</code>	Returns the values of the dictionary in a <code>dict_values</code> object
<code>items</code>	<code>adict.items()</code>	Returns the key-value pairs in a <code>dict_items</code> object
<code>get</code>	<code>adict.get(k)</code>	Returns the value associated with <code>k</code> , <code>None</code> otherwise
<code>get</code>	<code>adict.get(k,alt)</code>	Returns the value associated with <code>k</code> , <code>alt</code> otherwise

表格 8 Methods Provided by Dictionaries in Python



代码 6 (scratch\_01\_01)

### 1.7.2. 输入与输出

We often have a need to interact with users, either to get data or to provide some sort of result.

Most programs today use a dialog box as a way of asking the user to provide some type of input. While Python does have a way to create dialog boxes, there is a much simpler function that we can use. Python provides us with a function that allows us to ask a user to enter some data and returns a reference to the data in the form of a string. The function is called `input`.

Python's `input` function takes a single parameter that is a string. This string is often called the **prompt** because it contains some helpful text prompting the user to enter something. For example, you might call `input` as follows:

```
aName = input('Please enter your name: ')
```

Now whatever the user types after the prompt will be stored in the `aName` variable. Using the `input` function, we can easily write instructions that will prompt the user to enter data and then incorporate that data into further processing. For example, in the following two statements, the first asks the user for their name and the second prints the result of some simple processing based on the string that is provided.

```
aName = input("Please enter your name ")
print("Your name in all capitals is", aName.upper(),
      "and has length", len(aName))
```

代码 7 The `input` Function Returns a String (strstuff)

It is important to note that the value returned from the `input` function will be a string representing the exact characters that were entered after the prompt. If you want this string interpreted as another type, you must provide the type conversion explicitly. In the statements below, the string that is entered by the user is converted to a float so that it can be used in further arithmetic processing.

```
sradius = input("Please enter the radius of the circle ")
radius = float(sradius)
diameter = 2 * radius
```

### 1.7.2.1. 字符串格式化输出

We have already seen that the `print` function provides a very simple way to output values from a Python program. `print` takes zero or more parameters and displays them using a single blank as the default separator. It is possible to change the separator character by setting the `sep` argument. In addition, each `print` ends with a newline character by default. This behavior can be changed by setting the `end` argument. These variations are shown in the following session:

```
>>> print("Hello")
Hello
>>> print("Hello","World")
Hello World
>>> print("Hello","World", sep="****")
Hello***World
>>> print("Hello","World", end="****")
Hello World***>>>
```

It is often useful to have more control over the look of your output. Fortunately, Python provides us with an alternative called **formatted strings**. A formatted string is a template in which words or spaces that will remain constant are combined with placeholders for variables that will be inserted into the string. For example, the statement

```
print(aName, "is", age, "years old.")
```

contains the words `is` and `years old`, but the name and the age will change depending on the variable values at the time of execution. Using a formatted string, we write the previous statement as

```
print("%s is %d years old." % (aName, age))
```

This simple example illustrates a new string expression. The `%` operator is a string operator called the **format operator**. The left side of the expression holds the template or format string, and the right side holds a collection of values that will be substituted into the format string. Note that the number of values in the collection on the right side corresponds with the number of `%` characters in the format string. Values are taken—in order, left to right—from the collection and inserted into the format string.

Let's look at both sides of this formatting expression in more detail. The format string may contain one or more conversion specifications. A conversion character tells the format operator what type of value is going to be inserted into that position in the string. In the example above, the `%s` specifies a string, while the `%d` specifies an integer. Other possible type specifications include `i`, `u`, `f`, `e`, `g`, `c`, or `%`. Table 9 summarizes all of the various type specifications.

Character	Output Format
<code>d</code> , <code>i</code>	Integer
<code>u</code>	Unsigned integer
<code>f</code>	Floating point as m.ddddd
<code>e</code>	Floating point as m.dddddE+/-xx
<code>E</code>	Floating point as m.dddddE+/-xx
<code>g</code>	Use <code>%e</code> for exponents less than <code>-4</code> or greater than <code>+5</code> , otherwise use <code>%f</code>



<code>c</code>	Single character
<code>s</code>	String, or any Python data object that can be converted to a string by using the <code>str</code> function.
<code>%</code>	Insert a literal <code>%</code> character

表格 9 String Formatting Conversion Characters

In addition to the format character, you can also include a format modifier between the `%` and the format character. Format modifiers may be used to left-justify or right-justify the value with a specified field width. Modifiers can also be used to specify the field width along with a number of digits after the decimal point. Table 10 explains these format modifiers

Modifier	Example	Description
number	<code>%20d</code>	Put the value in a field width of 20
<code>-</code>	<code>%-20d</code>	Put the value in a field 20 characters wide, left-justified
<code>+</code>	<code>%+20d</code>	Put the value in a field 20 characters wide, right-justified
<code>0</code>	<code>%020d</code>	Put the value in a field 20 characters wide, fill in with leading zeros.
<code>.</code>	<code>%20.2f</code>	Put the value in a field 20 characters wide with 2 characters to the right of the decimal point.
<code>(name)</code>	<code>%(name)d</code>	Get the value from the supplied dictionary using <code>name</code> as the key.

表格 10 Additional formatting options

The right side of the format operator is a collection of values that will be inserted into the format string. The collection will be either a tuple or a dictionary. If the collection is a tuple, the values are inserted in order of position. That is, the first element in the tuple corresponds to the first format character in the format string. If the collection is a dictionary, the values are inserted according to their keys. In this case all format characters must use the `(name)` modifier to specify the name of the key.

```
>>> price = 24
>>> item = "banana"
>>> print("The %s costs %d cents"%(item,price))
The banana costs 24 cents
>>> print("The %+10s costs %5.2f cents"%(item,price))
The      banana costs 24.00 cents
>>> print("The %+10s costs %10.2f cents"%(item,price))
The      banana costs      24.00 cents
>>> itemdict = {"item":"banana","cost":24}
>>> print("The %(item)s costs %(cost)7.1f cents"%itemdict)
The banana costs      24.0 cents
>>>
```

In addition to format strings that use format characters and format modifiers, Python strings also include a `format` method that can be used in conjunction with a new `Formatter` class to

implement complex string formatting. More about these features can be found in the Python library reference manual.

### 1.7.3. 控制结构

As we noted earlier, algorithms require two important control structures: iteration and selection. Both of these are supported by Python in various forms. The programmer can choose the statement that is most useful for the given circumstance.

For iteration, Python provides a standard `while` statement and a very powerful `for` statement. The `while` statement repeats a body of code as long as a condition is true. For example,

```
>>> counter = 1
>>> while counter <= 5:
...     print("Hello, world")
...     counter = counter + 1
```

```
Hello, world
Hello, world
Hello, world
Hello, world
Hello, world
```

prints out the phrase “Hello, world” five times. The condition on the `while` statement is evaluated at the start of each repetition. If the condition is `True`, the body of the statement will execute. It is easy to see the structure of a Python `while` statement due to the mandatory indentation pattern that the language enforces.

The `while` statement is a very general purpose iterative structure that we will use in a number of different algorithms. In many cases, a compound condition will control the iteration. A fragment such as

```
while counter <= 10 and not done:
...
```

would cause the body of the statement to be executed only in the case where both parts of the condition are satisfied. The value of the variable `counter` would need to be less than or equal to 10 and the value of the variable `done` would need to be `False` (not `False` is `True`) so that `True` and `True` results in `True`.

Even though this type of construct is very useful in a wide variety of situations, another iterative structure, the `for` statement, can be used in conjunction with many of the Python collections. The `for` statement can be used to iterate over the members of a collection, so long as the collection is a sequence. So, for example,

```
>>> for item in [1,3,6,2,5]:
...     print(item)
...
1
3
6
2
5
```

assigns the variable `item` to be each successive value in the list `[1,3,6,2,5]`. The body of the iteration is then executed. This works for any collection that is a sequence (lists, tuples, and strings).

A common use of the `for` statement is to implement definite iteration over a range of values. The statement

```
>>> for item in range(5):
...     print(item**2)
...
0
1
4
9
16
>>>
```

will perform the `print` function five times. The `range` function will return a range object representing the sequence 0,1,2,3,4 and each value will be assigned to the variable `item`. This value is then squared and printed.

The other very useful version of this iteration structure is used to process each character of a string. The following code fragment iterates over a list of strings and for each string processes each character by appending it to a list. The result is a list of all the letters in all of the words.

```
wordlist = ['cat','dog','rabbit']
letterlist = [ ]
for aword in wordlist:
    for aletter in aword:
```

```
letterlist.append(aletter)
print(letterlist)
```

#### 代码 8 Processing Each Character in a List of Strings (intro\_8)

Selection statements allow programmers to ask questions and then, based on the result, perform different actions. Most programming languages provide two versions of this useful construct: the `ifelse` and the `if`. A simple example of a binary selection uses the `ifelse` statement.

```
if n<0:
    print("Sorry, value is negative")
else:
    print(math.sqrt(n))
```

In this example, the object referred to by `n` is checked to see if it is less than zero. If it is, a message is printed stating that it is negative. If it is not, the statement performs the `else` clause and computes the square root.

Selection constructs, as with any control construct, can be nested so that the result of one question helps decide whether to ask the next. For example, assume that `score` is a variable holding a reference to a score for a computer science test.

```
if score >= 90:
    print('A')
else:
    if score >=80:
        print('B')
    else:
        if score >= 70:
            print('C')
        else:
            if score >= 60:
                print('D')
            else:
                print('F')
```

This fragment will classify a value called `score` by printing the letter grade earned. If the score is greater than or equal to 90, the statement will print `A`. If it is not (`else`), the next question is asked. If the score is greater than or equal to 80 then it must be between 80 and 89 since the answer to the first question was false. In this case print `B` is printed. You can see that the Python indentation pattern helps to make sense of the association between `if` and `else` without requiring any additional syntactic elements.

An alternative syntax for this type of nested selection uses the `elif` keyword. The `else` and the next `if` are combined so as to eliminate the need for additional nesting levels. Note that the final `else` is still necessary to provide the default case if all other conditions fail.

```
if score >= 90:
    print('A')
elif score >=80:
    print('B')
elif score >= 70:
    print('C')
elif score >= 60:
    print('D')
else:
    print('F')
```

Python also has a single way selection construct, the `if` statement. With this statement, if the condition is true, an action is performed. In the case where the condition is false, processing simply continues on to the next statement after the `if`. For example, the following fragment will first check to see if the value of a variable `n` is negative. If it is, then it is modified by the absolute value function. Regardless, the next action is to compute the square root.

```
if n<0:
    n = abs(n)
print(math.sqrt(n))
```

### Self Check

Test your understanding of what we have covered so far by trying the following exercise. Modify the code from Activecode 8 so that the final list only contains a single copy of each letter.

```
# the answer is: ['c', 'a', 't', 'd', 'o', 'g', 'r', 'b', 'i']
```

代码 9 (self\_check\_1)

Returning to lists, there is an alternative method for creating a list that uses iteration and selection constructs. This is known as a list comprehension. A list comprehension allows you to easily create a list based on some processing or selection criteria. For example, if we would like to create a list of the first 10 perfect squares, we could use a `for` statement:

```
>>> sqlist=[]
>>> for x in range(1,11):
    sqlist.append(x*x)
```

```
>>> sqliist
[1, 4, 9, 16, 25, 36, 49, 64, 81, 100]
>>>
```

Using a list comprehension, we can do this in one step as

```
>>> sqliist=[x*x for x in range(1,11)]
>>> sqliist
[1, 4, 9, 16, 25, 36, 49, 64, 81, 100]
>>>
```

The variable `x` takes on the values 1 through 10 as specified by the `for` construct. The value of `x*x` is then computed and added to the list that is being constructed. The general syntax for a list comprehension also allows a selection criteria to be added so that only certain items get added. For example,

```
>>> sqliist=[x*x for x in range(1,11) if x%2 != 0]
>>> sqliist
[1, 9, 25, 49, 81]
>>>
```

This list comprehension constructed a list that only contained the squares of the odd numbers in the range from 1 to 10. Any sequence that supports iteration can be used within a list comprehension to construct a new list.

```
>>>[ch.upper() for ch in 'comprehension' if ch not in 'aeiou']
['C', 'M', 'P', 'R', 'H', 'N', 'S', 'N']
>>>
```

### Self Check

Test your understanding of list comprehensions by redoing Activecode 8 using list comprehensions. For an extra challenge, see if you can figure out how to remove the duplicates.

```
# the answer is: ['c', 'a', 't', 'd', 'o', 'g', 'r', 'a', 'b', 'b', 'i', 't']
```

### 1.7.4. 异常处理

There are two types of errors that typically occur when writing programs. The first, known as a syntax error, simply means that the programmer has made a mistake in the structure of a statement or expression. For example, it is incorrect to write a for statement and forget the colon.

```
>>> for i in range(10)
SyntaxError: invalid syntax (<pysHELL#61>, line 1)
```

In this case, the Python interpreter has found that it cannot complete the processing of this instruction since it does not conform to the rules of the language. Syntax errors are usually more frequent when you are first learning a language.

The other type of error, known as a logic error, denotes a situation where the program executes but gives the wrong result. This can be due to an error in the underlying algorithm or an error in your translation of that algorithm. In some cases, logic errors lead to very bad situations such as trying to divide by zero or trying to access an item in a list where the index of the item is outside the bounds of the list. In this case, the logic error leads to a runtime error that causes the program to terminate. These types of runtime errors are typically called **exceptions**.

Most of the time, beginning programmers simply think of exceptions as fatal runtime errors that cause the end of execution. However, most programming languages provide a way to deal with these errors that will allow the programmer to have some type of intervention if they so choose. In addition, programmers can create their own exceptions if they detect a situation in the program execution that warrants it.

When an exception occurs, we say that it has been “raised.” You can “handle” the exception that has been raised by using a **try** statement. For example, consider the following session that asks the user for an integer and then calls the square root function from the math library. If the user enters a value that is greater than or equal to 0, the print will show the square root. However, if the user enters a negative value, the square root function will report a **ValueError** exception.

```
>>> anumber = int(input("Please enter an integer "))
Please enter an integer -23
>>> print(math.sqrt(anumber))
Traceback (most recent call last):
  File "<pysHELL#102>", line 1, in <module>
    print(math.sqrt(anumber))
ValueError: math domain error
>>>
```

We can handle this exception by calling the print function from within a **try** block. A

corresponding `except` block “catches” the exception and prints a message back to the user in the event that an exception occurs. For example:

```
>>> try:
    print(math.sqrt(anumber))
except:
    print("Bad Value for square root")
    print("Using absolute value instead")
    print(math.sqrt(abs(anumber)))

Bad Value for square root
Using absolute value instead
4.79583152331
>>>
```

will catch the fact that an exception is raised by `sqrt` and will instead print the messages back to the user and use the absolute value to be sure that we are taking the square root of a non-negative number. This means that the program will not terminate but instead will continue on to the next statements.

It is also possible for a programmer to cause a runtime exception by using the `raise` statement. For example, instead of calling the square root function with a negative number, we could have checked the value first and then raised our own exception. The code fragment below shows the result of creating a new `RuntimeError` exception. Note that the program would still terminate but now the exception that caused the termination is something explicitly created by the programmer.

```
>>> if anumber < 0:
...     raise RuntimeError("You can't use a negative number")
... else:
...     print(math.sqrt(anumber))
...

Traceback (most recent call last):
  File "<stdin>", line 2, in <module>
RuntimeError: You can't use a negative number
>>>
```

There are many kinds of exceptions that can be raised in addition to the `RuntimeError` shown above. See the Python reference manual for a list of all the available exception types and for how to create your own.



### 1.7.5. 定义函数

The earlier example of procedural abstraction called upon a Python function called `sqrt` from the `math` module to compute the square root. In general, we can hide the details of any computation by defining a function. A function definition requires a name, a group of parameters, and a body. It may also explicitly return a value. For example, the simple function defined below returns the square of the value you pass into it.

```
>>> def square(n):  
...     return n**2  
...  
>>> square(3)  
9  
>>> square(square(3))  
81  
>>>
```

The syntax for this function definition includes the name, `square`, and a parenthesized list of formal parameters. For this function, `n` is the only formal parameter, which suggests that `square` needs only one piece of data to do its work. The details, hidden “inside the box,” simply compute the result of `n**2` and return it. We can invoke or call the `square` function by asking the Python environment to evaluate it, passing an actual parameter value, in this case, `3`. Note that the call to `square` returns an integer that can in turn be passed to another invocation.

We could implement our own square root function by using a well-known technique called “Newton’s Method.” Newton’s Method for approximating square roots performs an iterative computation that converges on the correct value. The equation  $newguess = 1/2 * (oldguess + n/oldguess)$  takes a value `n` and repeatedly guesses the square root by making each `newguess` the `oldguess` in the subsequent iteration. The initial guess used here is `n/2`. Listing 1 shows a function definition that accepts a value `n` and returns the square root of `n` after making 20 guesses. Again, the details of Newton’s Method are hidden inside the function definition and the user does not have to know anything about the implementation to use the function for its intended purpose. Listing 1 also shows the use of the `#` character as a comment marker. Any characters that follow the `#` on a line are ignored.

#### Listing 1

```
def squareroot(n):  
    root = n/2    #initial guess will be 1/2 of n  
    for k in range(20):  
        root = (1/2)*(root + (n / root))
```

```
return root
```

```
>>>squareroot(9)
3.0
>>>squareroot(4563)
67.549981495186216
>>>
```

### Self Check

Here's a self check that really covers everything so far. You may have heard of the infinite monkey theorem? The theorem states that a monkey hitting keys at random on a typewriter keyboard for an infinite amount of time will almost surely type a given text, such as the complete works of William Shakespeare. Well, suppose we replace a monkey with a Python function. How long do you think it would take for a Python function to generate just one sentence of Shakespeare? The sentence we'll shoot for is: "methinks it is like a weasel"

You're not going to want to run this one in the browser, so fire up your favorite Python IDE. The way we'll simulate this is to write a function that generates a string that is 27 characters long by choosing random letters from the 26 letters in the alphabet plus the space. We'll write another function that will score each generated string by comparing the randomly generated string to the goal.

A third function will repeatedly call generate and score, then if 100% of the letters are correct we are done. If the letters are not correct then we will generate a whole new string. To make it easier to follow your program's progress this third function should print out the best string generated so far and its score every 1000 tries.

### Self Check Challenge

See if you can improve upon the program in the self check by keeping letters that are correct and only modifying one character in the best string so far. This is a type of algorithm in the class of 'hill climbing' algorithms, that is we only keep the result if it is better than the previous one.

## 1.7.6. Python 面向对象编程：定义类

We stated earlier that Python is an object-oriented programming language. So far, we have used a number of built-in classes to show examples of data and control structures. One of the most powerful features in an object-oriented programming language is the ability to allow a programmer (problem solver) to create new classes that model data that is needed to solve the problem.

Remember that we use abstract data types to provide the logical description of what a data object looks like (its state) and what it can do (its methods). By building a class that implements an

abstract data type, a programmer can take advantage of the abstraction process and at the same time provide the details necessary to actually use the abstraction in a program. Whenever we want to implement an abstract data type, we will do so with a new class.

### 1.7.6.1. 示例：Fraction 类

A very common example to show the details of implementing a user-defined class is to construct a class to implement the abstract data type `Fraction`. We have already seen that Python provides a number of numeric classes for our use. There are times, however, that it would be most appropriate to be able to create data objects that “look like” fractions.

A fraction such as  $3/5$  consists of two parts. The top value, known as the numerator, can be any integer. The bottom value, called the denominator, can be any integer greater than 0 (negative fractions have a negative numerator). Although it is possible to create a floating point approximation for any fraction, in this case we would like to represent the fraction as an exact value.

The operations for the `Fraction` type will allow a `Fraction` data object to behave like any other numeric value. We need to be able to add, subtract, multiply, and divide fractions. We also want to be able to show fractions using the standard “slash” form, for example  $3/5$ . In addition, all fraction methods should return results in their lowest terms so that no matter what computation is performed, we always end up with the most common form.

In Python, we define a new class by providing a name and a set of method definitions that are syntactically similar to function definitions. For this example,

```
class Fraction:

    #the methods go here
```

provides the framework for us to define the methods. The first method that all classes should provide is the constructor. The constructor defines the way in which data objects are created. To create a `Fraction` object, we will need to provide two pieces of data, the numerator and the denominator. In Python, the constructor method is always called `__init__` (two underscores before and after `init`) and is shown in Listing 2.

#### Listing 2

```
class Fraction:

    def __init__(self,top,bottom):
```

```
self.num = top
self.den = bottom
```

Notice that the formal parameter list contains three items (`self`, `top`, `bottom`). `self` is a special parameter that will always be used as a reference back to the object itself. It must always be the first formal parameter; however, it will never be given an actual parameter value upon invocation. As described earlier, fractions require two pieces of state data, the numerator and the denominator. The notation `self.num` in the constructor defines the `fraction` object to have an internal data object called `num` as part of its state. Likewise, `self.den` creates the denominator. The values of the two formal parameters are initially assigned to the state, allowing the new `fraction` object to know its starting value.

To create an instance of the `Fraction` class, we must invoke the constructor. This happens by using the name of the class and passing actual values for the necessary state (note that we never directly invoke `__init__`). For example,

```
myfraction = Fraction(3,5)
```

creates an object called `myfraction` representing the fraction 3/5 (three-fifths). Figure 5 shows this object as it is now implemented.

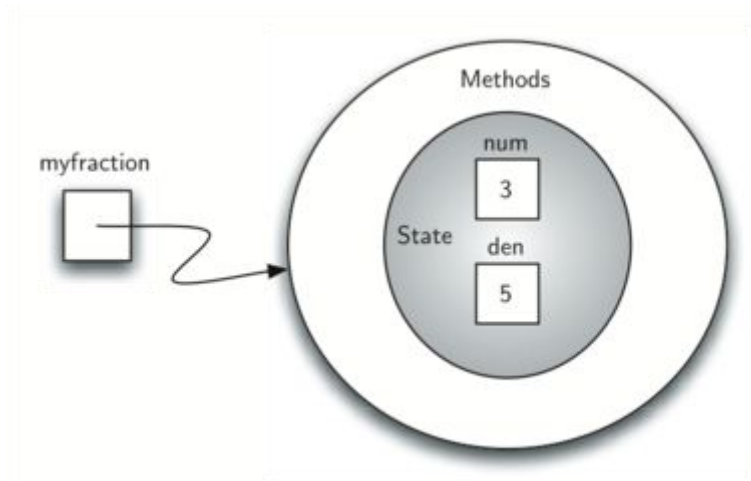


图 5 An Instance of the Fraction Class

The next thing we need to do is implement the behavior that the abstract data type requires. To begin, consider what happens when we try to print a `Fraction` object.

```
>>> myf = Fraction(3,5)
>>> print(myf)
<__main__.Fraction instance at 0x409b1acc>
```

The fraction object, `myf`, does not know how to respond to this request to print. The print function

requires that the object convert itself into a string so that the string can be written to the output. The only choice `myf` has is to show the actual reference that is stored in the variable (the address itself). This is not what we want.

There are two ways we can solve this problem. One is to define a method called `show` that will allow the `Fraction` object to print itself as a string. We can implement this method as shown in Listing 3. If we create a `Fraction` object as before, we can ask it to show itself, in other words, print itself in the proper format. Unfortunately, this does not work in general. In order to make printing work properly, we need to tell the `Fraction` class how to convert itself into a string. This is what the `print` function needs in order to do its job.

### Listing 3

```
def show(self):
    print(self.num,"/",self.den)
>>> myf = Fraction(3,5)
>>> myf.show()
3 / 5
>>> print(myf)
<__main__.Fraction instance at 0x40bce9ac>
>>>
```

In Python, all classes have a set of standard methods that are provided but may not work properly. One of these, `__str__`, is the method to convert an object into a string. The default implementation for this method is to return the instance address string as we have already seen. What we need to do is provide a “better” implementation for this method. We will say that this implementation overrides the previous one, or that it redefines the method’s behavior.

To do this, we simply define a method with the name `__str__` and give it a new implementation as shown in Listing 4. This definition does not need any other information except the special parameter `self`. In turn, the method will build a string representation by converting each piece of internal state data to a string and then placing a `/` character in between the strings using string concatenation. The resulting string will be returned any time a `Fraction` object is asked to convert itself to a string. Notice the various ways that this function is used.

### Listing 4

```
def __str__(self):
    return str(self.num)+"/"+str(self.den)
>>> myf = Fraction(3,5)
>>> print(myf)
3/5
>>> print("I ate", myf, "of the pizza")
```

```

I ate 3/5 of the pizza
>>> myf.__str__()
'3/5'
>>> str(myf)
'3/5'
>>>

```

We can override many other methods for our new Fraction class. Some of the most important of these are the basic arithmetic operations. We would like to be able to create two Fraction objects and then add them together using the standard “+” notation. At this point, if we try to add two fractions, we get the following:

```

>>> f1 = Fraction(1,4)
>>> f2 = Fraction(1,2)
>>> f1+f2

Traceback (most recent call last):
  File "<pyshell#173>", line 1, in -toplevel-
    f1+f2
TypeError: unsupported operand type(s) for +:
      'instance' and 'instance'
>>>

```

If you look closely at the error, you see that the problem is that the “+” operator does not understand the Fraction operands.

We can fix this by providing the Fraction class with a method that overrides the addition method. In Python, this method is called `__add__` and it requires two parameters. The first, `self`, is always needed, and the second represents the other operand in the expression. For example,

```

f1.__add__(f2)

```

would ask the Fraction object `f1` to add the Fraction object `f2` to itself. This can be written in the standard notation, `f1+f2`.

Two fractions must have the same denominator to be added. The easiest way to make sure they have the same denominator is to simply use the product of the two denominators as a common denominator so that  $a/b + c/d = ad/bd + cb/bd = (ad+cb)/bd$ . The implementation is shown in Listing 5. The addition function returns a new Fraction object with the numerator and denominator of the sum. We can use this method by writing a standard arithmetic expression involving fractions, assigning the result of the addition, and then printing our result.

### Listing 5

```
def __add__(self, otherfraction):  
  
    newnum = self.num*otherfraction.den + self.den*otherfraction.num  
    newden = self.den * otherfraction.den  
  
    return Fraction(newnum, newden)
```

```
>>> f1=Fraction(1,4)  
>>> f2=Fraction(1,2)  
>>> f3=f1+f2  
>>> print(f3)  
6/8  
>>>
```

The addition method works as we desire, but one thing could be better. Note that 6/8 is the correct result ( $1/4 + 1/2$ ) but that it is not in the “lowest terms” representation. The best representation would be 3/4. In order to be sure that our results are always in the lowest terms, we need a helper function that knows how to reduce fractions. This function will need to look for the greatest common divisor, or GCD. We can then divide the numerator and the denominator by the GCD and the result will be reduced to lowest terms.

The best-known algorithm for finding a greatest common divisor is Euclid’s Algorithm, which will be discussed in detail in Chapter 8. Euclid’s Algorithm states that the greatest common divisor of two integers  $m$  and  $n$  is  $n$  if  $n$  divides  $m$  evenly. However, if  $n$  does not divide  $m$  evenly, then the answer is the greatest common divisor of  $n$  and the remainder of  $m$  divided by  $n$ . We will simply provide an iterative implementation here (see ActiveCode 1). Note that this implementation of the GCD algorithm only works when the denominator is positive. This is acceptable for our fraction class because we have said that a negative fraction will be represented by a negative numerator.

```
def gcd(m,n):  
    while m%n != 0:  
        oldm = m  
        oldn = n  
  
        m = oldn  
        n = oldm%oldn  
    return n  
  
print gcd(20,10)
```

## 代码 11 The Greatest Common Divisor Function (gcd\_cl)

Now we can use this function to help reduce any fraction. To put a fraction in lowest terms, we will divide the numerator and the denominator by their greatest common divisor. So, for the fraction 6/8, the greatest common divisor is 2. Dividing the top and the bottom by 2 creates a new fraction, 3/4 (see Listing 6).

### Listing 6

```
def __add__(self, otherfraction):
    newnum = self.num*otherfraction.den + self.den*otherfraction.num
    newden = self.den * otherfraction.den
    common = gcd(newnum,newden)
    return Fraction(newnum//common,newden//common)
```

```
>>> f1=Fraction(1,4)
>>> f2=Fraction(1,2)
>>> f3=f1+f2
>>> print(f3)
3/4
>>>
```

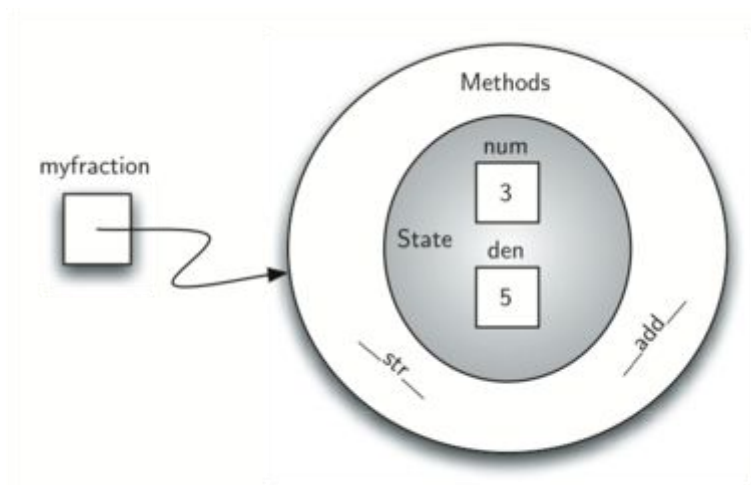


图 6 An Instance of the Fraction Class with Two Methods

Our Fraction object now has two very useful methods and looks like Figure 6. An additional group of methods that we need to include in our example Fraction class will allow two fractions to compare themselves to one another. Assume we have two Fraction objects, f1 and f2. f1==f2 will only be True if they are references to the same object. Two different objects with the same numerators and denominators would not be equal under this implementation. This is called shallow equality (see Figure 7).



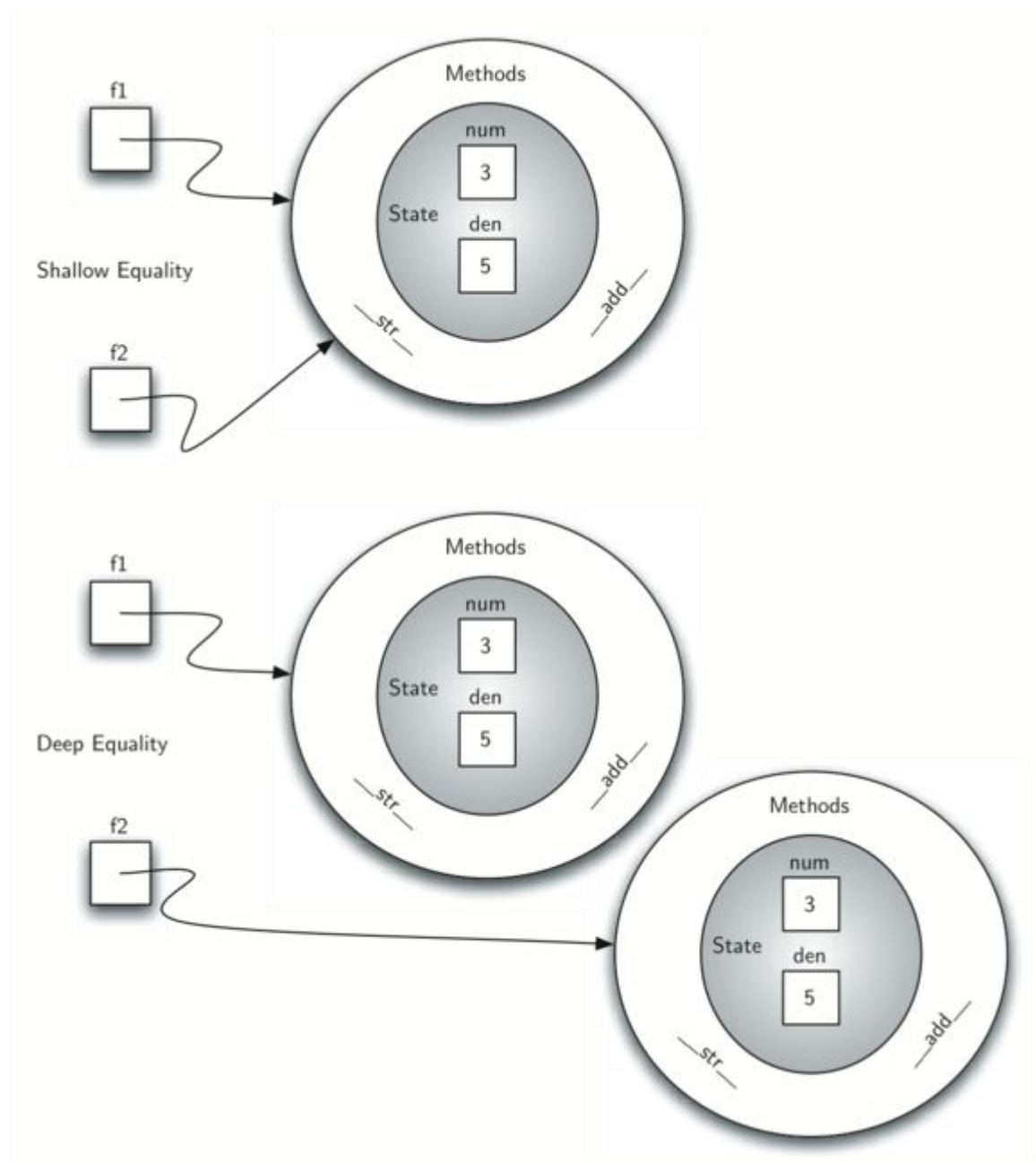


图 7 Shallow Equality Versus Deep Equality

We can create deep equality (see Figure 7) – equality by the same value, not the same reference – by overriding the `__eq__` method. The `__eq__` method is another standard method available in any class. The `__eq__` method compares two objects and returns `True` if their values are the same, `False` otherwise.

In the `Fraction` class, we can implement the `__eq__` method by again putting the two fractions in common terms and then comparing the numerators (see Listing 7). It is important to note that there are other relational operators that can be overridden. For example, the `__le__` method provides the less than or equal functionality.

### Listing 7

```
def __eq__(self, other):  
    firstnum = self.num * other.den  
    secondnum = other.num * self.den  
  
    return firstnum == secondnum
```

The complete Fraction class, up to this point, is shown in ActiveCode 2. We leave the remaining arithmetic and relational methods as exercises.

### Self Check

To make sure you understand how operators are implemented in Python classes, and how to properly write methods, write some methods to implement `*`, `/`, and `-`. Also implement comparison operators `>` and `<`

代码 12 (self\_check\_4)

### 1.7.6.2. 继承：逻辑门与门电路

Our final section will introduce another important aspect of object-oriented programming. Inheritance is the ability for one class to be related to another class in much the same way that people can be related to one another. Children inherit characteristics from their parents. Similarly, Python child classes can inherit characteristic data and behavior from a parent class. These classes are often referred to as **subclasses** and **superclasses**.

Figure 8 shows the built-in Python collections and their relationships to one another. We call a relationship structure such as this an inheritance hierarchy. For example, the list is a child of the sequential collection. In this case, we call the list the child and the sequence the parent (or subclass list and superclass sequence). This is often referred to as an IS-A Relationship (the list IS-A sequential collection). This implies that lists inherit important characteristics from sequences, namely the ordering of the underlying data and operations such as concatenation, repetition, and indexing.

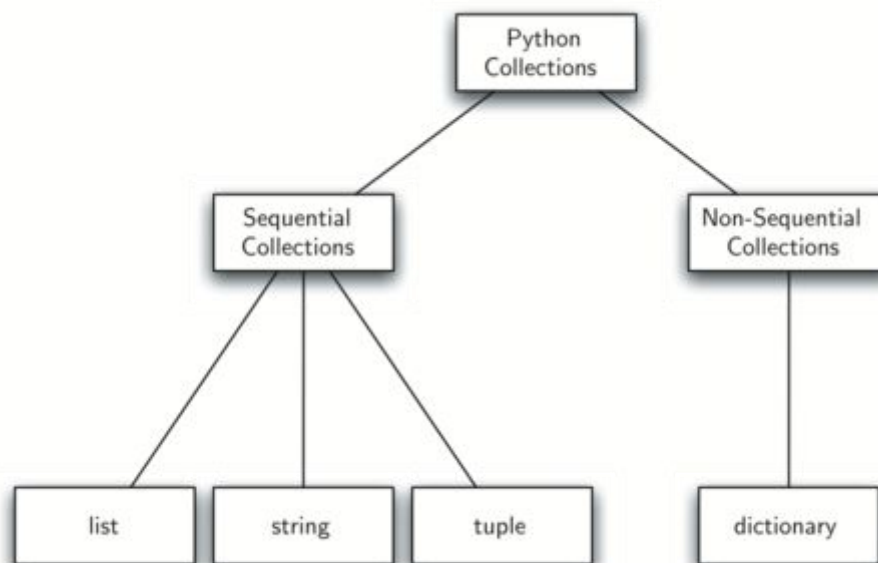


图 8 An Inheritance Hierarchy for Python Collections

Lists, tuples, and strings are all types of sequential collections. They all inherit common data organization and operations. However, each of them is distinct based on whether the data is homogeneous and whether the collection is immutable. The children all gain from their parents but distinguish themselves by adding additional characteristics.

By organizing classes in this hierarchical fashion, object-oriented programming languages allow previously written code to be extended to meet the needs of a new situation. In addition, by organizing data in this hierarchical manner, we can better understand the relationships that exist. We can be more efficient in building our abstract representations.

To explore this idea further, we will construct a simulation, an application to simulate digital circuits. The basic building block for this simulation will be the logic gate. These electronic switches represent boolean algebra relationships between their input and their output. In general, gates have a single output line. The value of the output is dependent on the values given on the input lines.

AND gates have two input lines, each of which can be either 0 or 1 (representing False or True, respectively). If both of the input lines have the value 1, the resulting output is 1. However, if either or both of the input lines is 0, the result is 0. OR gates also have two input lines and produce a 1 if one or both of the input values is a 1. In the case where both input lines are 0, the result is 0.

NOT gates differ from the other two gates in that they only have a single input line. The output value is simply the opposite of the input value. If 0 appears on the input, 1 is produced on the output. Similarly, 1 produces 0. Figure 9 shows how each of these gates is typically represented. Each gate also has a truth table of values showing the input-to-output mapping that is performed by the gate.

../\_images/truthtable.png

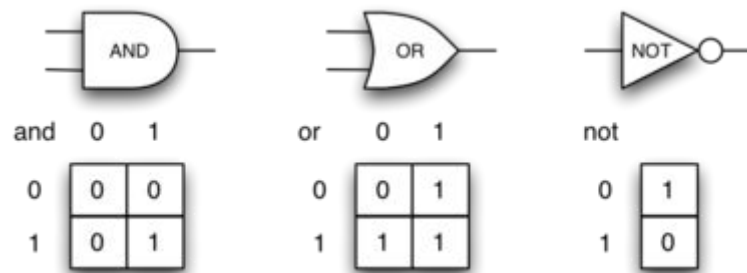


图 9 Three Types of Logic Gates

By combining these gates in various patterns and then applying a set of input values, we can build circuits that have logical functions. Figure 10 shows a circuit consisting of two AND gates, one OR gate, and a single NOT gate. The output lines from the two AND gates feed directly into the OR gate, and the resulting output from the OR gate is given to the NOT gate. If we apply a set of input values to the four input lines (two for each AND gate), the values are processed and a result appears at the output of the NOT gate. Figure 10 also shows an example with values.

../\_images/circuit1.png

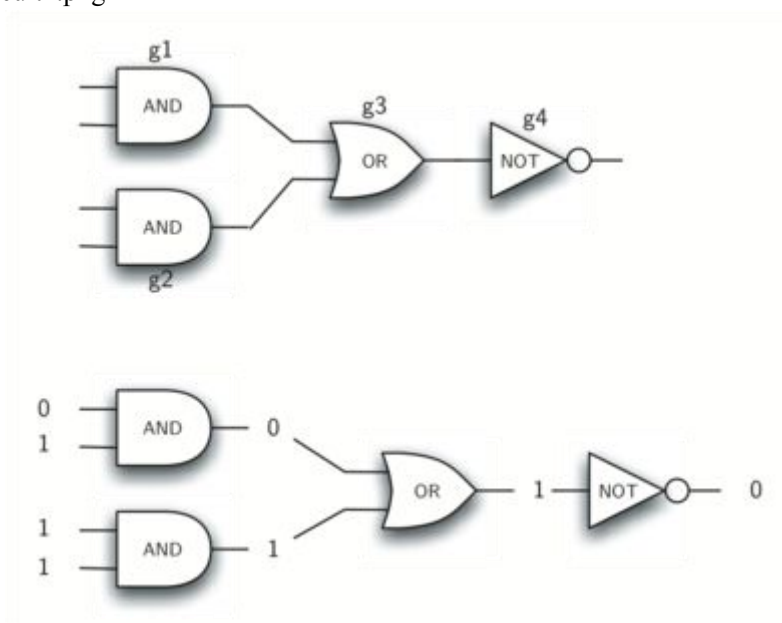


图 10 Circuit

In order to implement a circuit, we will first build a representation for logic gates. Logic gates are easily organized into a class inheritance hierarchy as shown in Figure 11. At the top of the hierarchy, the LogicGate class represents the most general characteristics of logic gates: namely, a label for the gate and an output line. The next level of subclasses breaks the logic gates into two families, those that have one input line and those that have two. Below that, the specific logic functions of each appear.

../\_images/gates.png

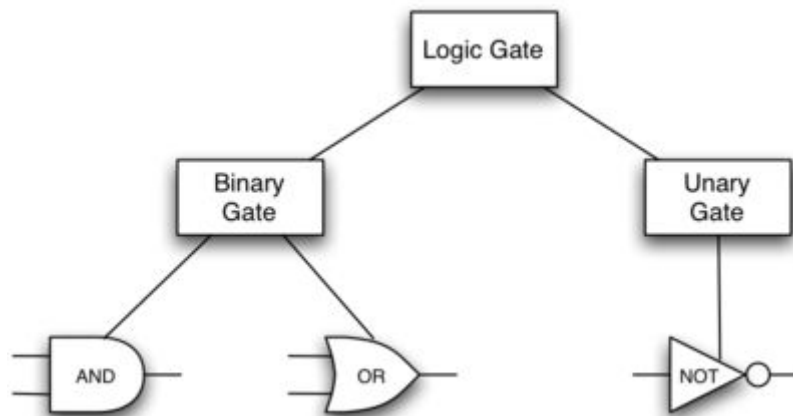


图 11 An Inheritance Hierarchy for Logic Gates

We can now start to implement the classes by starting with the most general, `LogicGate`. As noted earlier, each gate has a label for identification and a single output line. In addition, we need methods to allow a user of a gate to ask the gate for its label.

The other behavior that every logic gate needs is the ability to know its output value. This will require that the gate perform the appropriate logic based on the current input. In order to produce output, the gate needs to know specifically what that logic is. This means calling a method to perform the logic computation. The complete class is shown in Listing 8.

#### Listing 8

```
class LogicGate:

    def __init__(self,n):
        self.label = n
        self.output = None

    def getLabel(self):
        return self.label

    def getOutput(self):
        self.output = self.performGateLogic()
        return self.output
```

At this point, we will not implement the `performGateLogic` function. The reason for this is that we do not know how each gate will perform its own logic operation. Those details will be included by each individual gate that is added to the hierarchy. This is a very powerful idea in object-oriented programming. We are writing a method that will use code that does not exist yet. The parameter `self` is a reference to the actual gate object invoking the method. Any new logic gate that gets added to the hierarchy will simply need to implement the `performGateLogic` function and it will

be used at the appropriate time. Once done, the gate can provide its output value. This ability to extend a hierarchy that currently exists and provide the specific functions that the hierarchy needs to use the new class is extremely important for reusing existing code.

We categorized the logic gates based on the number of input lines. The AND gate has two input lines. The OR gate also has two input lines. NOT gates have one input line. The BinaryGate class will be a subclass of LogicGate and will add two input lines. The UnaryGate class will also subclass LogicGate but will have only a single input line. In computer circuit design, these lines are sometimes called “pins” so we will use that terminology in our implementation.

#### Listing 9

```
class BinaryGate(LogicGate):

    def __init__(self,n):
        LogicGate.__init__(self,n)

        self.pinA = None
        self.pinB = None

    def getPinA(self):
        return int(input("Enter Pin A input for gate "+
self.getLabel()+"-->"))

    def getPinB(self):
        return int(input("Enter Pin B input for gate "+
self.getLabel()+"-->"))
```

#### Listing 10

```
class UnaryGate(LogicGate):

    def __init__(self,n):
        LogicGate.__init__(self,n)

        self.pin = None

    def getPin(self):
        return int(input("Enter Pin input for gate "+
self.getLabel()+"-->"))
```

Listing 9 and Listing 10 implement these two classes. The constructors in both of these classes start with an explicit call to the constructor of the parent class using the parent's `__init__` method.

When creating an instance of the BinaryGate class, we first want to initialize any data items that are inherited from LogicGate. In this case, that means the label for the gate. The constructor then goes on to add the two input lines (pinA and pinB). This is a very common pattern that you should always use when building class hierarchies. Child class constructors need to call parent class constructors and then move on to their own distinguishing data.

Python also has a function called super which can be used in place of explicitly naming the parent class. This is a more general mechanism, and is widely used, especially when a class has more than one parent. But, this is not something we are going to discuss in this introduction. For example in our example above LogicGate.\_\_init\_\_(self,n) could be replaced with super(UnaryGate,self).\_\_init\_\_(n).

The only behavior that the BinaryGate class adds is the ability to get the values from the two input lines. Since these values come from some external place, we will simply ask the user via an input statement to provide them. The same implementation occurs for the UnaryGate class except that there is only one input line.

Now that we have a general class for gates depending on the number of input lines, we can build specific gates that have unique behavior. For example, the AndGate class will be a subclass of BinaryGate since AND gates have two input lines. As before, the first line of the constructor calls upon the parent class constructor (BinaryGate), which in turn calls its parent class constructor (LogicGate). Note that the AndGate class does not provide any new data since it inherits two input lines, one output line, and a label.

#### Listing 11

```
class AndGate(BinaryGate):

    def __init__(self,n):
        BinaryGate.__init__(self,n)

    def performGateLogic(self):

        a = self.getPinA()
        b = self.getPinB()
        if a==1 and b==1:
            return 1
        else:
            return 0
```

The only thing AndGate needs to add is the specific behavior that performs the boolean operation that was described earlier. This is the place where we can provide the performGateLogic method. For an AND gate, this method first must get the two input values and then only return 1 if both input values are 1. The complete class is shown in Listing 11.

We can show the AndGate class in action by creating an instance and asking it to compute its output. The following session shows an AndGate object, g1, that has an internal label "G1". When we invoke the getOutput method, the object must first call its performGateLogic method which in turn queries the two input lines. Once the values are provided, the correct output is shown.

```
>>> g1 = AndGate("G1")
>>> g1.getOutput()
Enter Pin A input for gate G1-->1
Enter Pin B input for gate G1-->0
0
```

The same development can be done for OR gates and NOT gates. The OrGate class will also be a subclass of BinaryGate and the NotGate class will extend the UnaryGate class. Both of these classes will need to provide their own performGateLogic functions, as this is their specific behavior.

We can use a single gate by first constructing an instance of one of the gate classes and then asking the gate for its output (which will in turn need inputs to be provided). For example:

```
>>> g2 = OrGate("G2")
>>> g2.getOutput()
Enter Pin A input for gate G2-->1
Enter Pin B input for gate G2-->1
1
>>> g2.getOutput()
Enter Pin A input for gate G2-->0
Enter Pin B input for gate G2-->0
0
>>> g3 = NotGate("G3")
>>> g3.getOutput()
Enter Pin input for gate G3-->0
1
```

Now that we have the basic gates working, we can turn our attention to building circuits. In order to create a circuit, we need to connect gates together, the output of one flowing into the input of another. To do this, we will implement a new class called Connector.

The Connector class will not reside in the gate hierarchy. It will, however, use the gate hierarchy in that each connector will have two gates, one on either end (see Figure 12). This relationship is very important in object-oriented programming. It is called the HAS-A Relationship. Recall earlier that we used the phrase “IS-A Relationship” to say that a child class is related to a parent class, for example UnaryGate IS-A LogicGate.



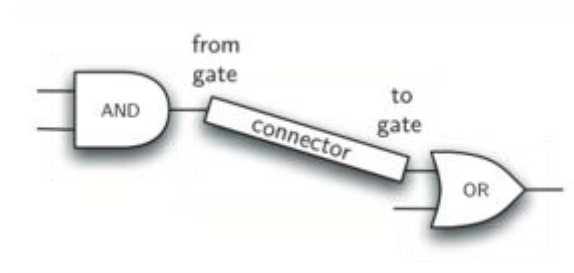


图 12 A Connector Connects the Output of One Gate to the Input of Another

Now, with the Connector class, we say that a Connector HAS-A LogicGate meaning that connectors will have instances of the LogicGate class within them but are not part of the hierarchy. When designing classes, it is very important to distinguish between those that have the IS-A relationship (which requires inheritance) and those that have HAS-A relationships (with no inheritance).

Listing 12 shows the Connector class. The two gate instances within each connector object will be referred to as the fromgate and the togate, recognizing that data values will “flow” from the output of one gate into an input line of the next. The call to setNextPin is very important for making connections (see Listing 13). We need to add this method to our gate classes so that each togate can choose the proper input line for the connection.

#### Listing 12

```
class Connector:

    def __init__(self, fgate, tgate):
        self.fromgate = fgate
        self.togate = tgate

        tgate.setNextPin(self)

    def getFrom(self):
        return self.fromgate

    def getTo(self):
        return self.togate
```

In the BinaryGate class, for gates with two possible input lines, the connector must be connected to only one line. If both of them are available, we will choose pinA by default. If pinA is already connected, then we will choose pinB. It is not possible to connect to a gate with no available input lines.

#### Listing 13

```
def setNextPin(self,source):
    if self.pinA == None:
        self.pinA = source
    else:
        if self.pinB == None:
            self.pinB = source
        else:
            raise RuntimeError("Error: NO EMPTY PINS")
```

Now it is possible to get input from two places: externally, as before, and from the output of a gate that is connected to that input line. This requires a change to the `getPinA` and `getPinB` methods (see Listing 14). If the input line is not connected to anything (`None`), then ask the user externally as before. However, if there is a connection, the connection is accessed and `fromgate`'s output value is retrieved. This in turn causes that gate to process its logic. This continues until all input is available and the final output value becomes the required input for the gate in question. In a sense, the circuit works backwards to find the input necessary to finally produce output.

#### Listing 14

```
def getPinA(self):
    if self.pinA == None:
        return input("Enter Pin A input for gate " + self.getName()+"-->")
    else:
        return self.pinA.getFrom().getOutput()
```

The following fragment constructs the circuit shown earlier in the section:

```
>>> g1 = AndGate("G1")
>>> g2 = AndGate("G2")
>>> g3 = OrGate("G3")
>>> g4 = NotGate("G4")
>>> c1 = Connector(g1,g3)
>>> c2 = Connector(g2,g3)
>>> c3 = Connector(g3,g4)
```

The outputs from the two AND gates (`g1` and `g2`) are connected to the OR gate (`g3`) and that output is connected to the NOT gate (`g4`). The output from the NOT gate is the output of the entire circuit. For example:

```
>>> g4.getOutput()
Pin A input for gate G1-->0
```

```
Pin B input for gate G1-->1
Pin A input for gate G2-->1
Pin B input for gate G2-->1
0
```

Try it yourself using ActiveCode 4.

```
class LogicGate:

    def __init__(self,n):
        self.name = n
        self.output = None

    def getName(self):
        return self.name

    def getOutput(self):
        self.output = self.performGateLogic()
        return self.output

class BinaryGate(LogicGate):

    def __init__(self,n):
        LogicGate.__init__(self,n)

        self.pinA = None
        self.pinB = None

    def getPinA(self):
        if self.pinA == None:
            return int(input("Enter Pin A input for gate\n"+self.getName()+"-->"))
        else:
            return self.pinA.getFrom().getOutput()

    def getPinB(self):
        if self.pinB == None:
            return int(input("Enter Pin B input for gate\n"+self.getName()+"-->"))
        else:
            return self.pinB.getFrom().getOutput()
```

```

def setNextPin(self,source):
    if self.pinA == None:
        self.pinA = source
    else:
        if self.pinB == None:
            self.pinB = source
        else:
            print("Cannot Connect: NO EMPTY PINS on this gate")

class AndGate(BinaryGate):

    def __init__(self,n):
        BinaryGate.__init__(self,n)

    def performGateLogic(self):

        a = self.getPinA()
        b = self.getPinB()
        if a==1 and b==1:
            return 1
        else:
            return 0

class OrGate(BinaryGate):

    def __init__(self,n):
        BinaryGate.__init__(self,n)

    def performGateLogic(self):

        a = self.getPinA()
        b = self.getPinB()
        if a ==1 or b==1:
            return 1
        else:
            return 0

class UnaryGate(LogicGate):

    def __init__(self,n):
        LogicGate.__init__(self,n)

```

```

        self.pin = None

    def getPin(self):
        if self.pin == None:
            return int(input("Enter Pin input for gate
"+self.getName()+"-->"))
        else:
            return self.pin.getFrom().getOutput()

    def setNextPin(self,source):
        if self.pin == None:
            self.pin = source
        else:
            print("Cannot Connect: NO EMPTY PINS on this gate")

class NotGate(UnaryGate):

    def __init__(self,n):
        UnaryGate.__init__(self,n)

    def performGateLogic(self):
        if self.getPin():
            return 0
        else:
            return 1

class Connector:

    def __init__(self, fgate, tgate):
        self.fromgate = fgate
        self.togate = tgate

        tgate.setNextPin(self)

    def getFrom(self):
        return self.fromgate

    def getTo(self):
        return self.togate

def main():

```

```

g1 = AndGate("G1")
g2 = AndGate("G2")
g3 = OrGate("G3")
g4 = NotGate("G4")
c1 = Connector(g1,g3)
c2 = Connector(g2,g3)
c3 = Connector(g3,g4)
print(g4.getOutput())

main()

```

代码 13 The Complete Circuit Program. (complete\_circuit)

### Self Check

Create a two new gate classes, one called NorGate the other called NandGate. NandGates work like AndGates that have a Not attached to the output. NorGates work like OrGates that have a Not attached to the output.

Create a series of gates that prove the following equality  $\text{NOT}((A \text{ and } B) \text{ or } (C \text{ and } D))$  is that same as  $\text{NOT}(A \text{ and } B) \text{ and } \text{NOT}(C \text{ and } D)$ . Make sure to use some of your new gates in the simulation.

代码 14 (self\_check\_5)

## 1.8. 小结

- 计算机科学是研究问题求解的学科；
- 计算机科学采用“抽象”作为表示过程与数据的工具；
- 采用“抽象数据类型”，程序员可以通过隐藏数据细节来控制问题域的复杂度；
- Python 是一个强大的、而又易于使用的面向对象程序设计语言；
- 列表、元组和串，是 Python 内置的有序集类型；
- 字典与集合，是数据的无序集类型；
- 采用“类”，程序员可以具体实现一个抽象数据类型；
- 程序员可以创建新方法，也可以重载已有的标准方法；
- 一个类构造器在处理自身的数据和行为之前，总会调用其父类的构造器。

## 1.9. 关键词

abstract data type	abstraction	algorithm
class	computable	data abstraction

data structure	data type	deep equality
dictionary	encapsulation	exception
format operator	formatted strings	HAS-A relationship
implementation-independent	information hiding	inheritance
inheritance hierarchy	interface	IS-A relationship
list	list comprehension	method
mutability	object	procedural abstraction
programming	prompt	self
shallow equality	simulation	string
subclass	superclass	truth table

## 1.10. 问题讨论

1. Construct a class hierarchy for people on a college campus. Include faculty, staff, and students. What do they have in common? What distinguishes them from one another?
2. Construct a class hierarchy for bank accounts.
3. Construct a class hierarchy for different types of computers.
4. Using the classes provided in the chapter, interactively construct a circuit and test it.

## 1.11. 编程练习

1. Implement the simple methods `getNum` and `getDen` that will return the numerator and denominator of a fraction.
2. In many ways it would be better if all fractions were maintained in lowest terms right from the start. Modify the constructor for the `Fraction` class so that GCD is used to reduce fractions immediately. Notice that this means the `__add__` function no longer needs to reduce. Make the necessary modifications.
3. Implement the remaining simple arithmetic operators (`__sub__`, `__mul__`, and `__truediv__`).
4. Implement the remaining relational operators (`__gt__`, `__ge__`, `__lt__`, `__le__`, and `__ne__`).
5. Modify the constructor for the fraction class so that it checks to make sure that the numerator and denominator are both integers. If either is not an integer the constructor should raise an exception.
6. In the definition of fractions we assumed that negative fractions have a negative numerator and a positive denominator. Using a negative denominator would cause some of the relational operators to give incorrect results. In general, this is an unnecessary constraint. Modify the constructor to allow the user to pass a negative denominator so that all of the operators continue to work properly.
7. Research the `__radd__` method. How does it differ from `__add__`? When is it used? Implement `__radd__`.
8. Repeat the last question but this time consider the `__iadd__` method.
9. Research the `__repr__` method. How does it differ from `__str__`? When is it used? Implement `__repr__`.
10. Research other types of gates that exist (such as NAND, NOR, and XOR). Add them to the

circuit hierarchy. How much additional coding did you need to do?

11. The most simple arithmetic circuit is known as the half-adder. Research the simple half-adder circuit. Implement this circuit.
12. Now extend that circuit and implement an 8 bit full-adder.
13. The circuit simulation shown in this chapter works in a backward direction. In other words, given a circuit, the output is produced by working back through the input values, which in turn cause other outputs to be queried. This continues until external input lines are found, at which point the user is asked for values. Modify the implementation so that the action is in the forward direction; upon receiving inputs the circuit produces an output.
14. Design a class to represent a playing card. Now design a class to represent a deck of cards. Using these two classes, implement a favorite card game.
15. Find a Sudoku puzzle in the local newspaper. Write a program to solve the puzzle.



## 2. 算法分析

### 2.1. 目标

- 了解为何算法分析非常重要；
- 能够采用“大 O”方法来描述算法执行时间；
- 了解在 Python 列表和字典类型中通用操作执行时间的“大 O”级别；
- 了解 Python 数据类型的具体实现对算法分析的影响；
- 了解如何对简单 Python 程序进行执行时间检测。

### 2.2. 什么是算法分析

It is very common for beginning computer science students to compare their programs with one another. You may also have noticed that it is common for computer programs to look very similar, especially the simple ones. An interesting question often arises. When two programs solve the same problem but look different, is one program better than the other?

In order to answer this question, we need to remember that there is an important difference between a program and the underlying algorithm that the program is representing. As we stated in Chapter 1, an algorithm is a generic, step-by-step list of instructions for solving a problem. It is a method for solving any instance of the problem such that given a particular input, the algorithm produces the desired result. A program, on the other hand, is an algorithm that has been encoded into some programming language. There may be many programs for the same algorithm, depending on the programmer and the programming language being used.

To explore this difference further, consider the function shown in ActiveCode 1. This function solves a familiar problem, computing the sum of the first  $n$  integers. The algorithm uses the idea of an accumulator variable that is initialized to 0. The solution then iterates through the  $n$  integers, adding each to the accumulator.

```
def sumOfN(n):  
    theSum = 0  
    for i in range(1,n+1):  
        theSum = theSum + i  
  
    return theSum  
  
print(sumOfN(10))
```

代码 15 Summation of the First  $n$  Integers (active1)

Now look at the function in ActiveCode 2. At first glance it may look strange, but upon further inspection you can see that this function is essentially doing the same thing as the previous one. The reason this is not obvious is poor coding. We did not use good identifier names to assist with readability, and we used an extra assignment statement during the accumulation step that was not really necessary.

```
def foo(tom):  
    fred = 0  
    for bill in range(1,tom+1):  
        barney = bill  
        fred = fred + barney  
  
    return fred  
  
print(foo(10))
```

代码 16 Another Summation of the First n Integers (active2)

The question we raised earlier asked whether one function is better than another. The answer depends on your criteria. The function `sumOfN` is certainly better than the function `foo` if you are concerned with readability. In fact, you have probably seen many examples of this in your introductory programming course since one of the goals there is to help you write programs that are easy to read and easy to understand. In this course, however, we are also interested in characterizing the algorithm itself. (We certainly hope that you will continue to strive to write readable, understandable code.)

Algorithm analysis is concerned with comparing algorithms based upon the amount of computing resources that each algorithm uses. We want to be able to consider two algorithms and say that one is better than the other because it is more efficient in its use of those resources or perhaps because it simply uses fewer. From this perspective, the two functions above seem very similar. They both use essentially the same algorithm to solve the summation problem.

At this point, it is important to think more about what we really mean by computing resources. There are two different ways to look at this. One way is to consider the amount of space or memory an algorithm requires to solve the problem. The amount of space required by a problem solution is typically dictated by the problem instance itself. Every so often, however, there are algorithms that have very specific space requirements, and in those cases we will be very careful to explain the variations.

As an alternative to space requirements, we can analyze and compare algorithms based on the amount of time they require to execute. This measure is sometimes referred to as the “execution time” or “running time” of the algorithm. One way we can measure the execution time for the function `sumOfN` is to do a benchmark analysis. This means that we will track the actual time required for the program to compute its result. In Python, we can benchmark a function by noting the starting time and ending time with respect to the system we are using. In the `time` module there

is a function called `time` that will return the current system clock time in seconds since some arbitrary starting point. By calling this function twice, at the beginning and at the end, and then computing the difference, we can get an exact number of seconds (fractions in most cases) for execution.

### Listing 1

```
import time

def sumOfN2(n):
    start = time.time()

    theSum = 0
    for i in range(1,n+1):
        theSum = theSum + i

    end = time.time()

    return theSum,end-start
```

Listing 1 shows the original `sumOfN` function with the timing calls embedded before and after the summation. The function returns a tuple consisting of the result and the amount of time (in seconds) required for the calculation. If we perform 5 invocations of the function, each computing the sum of the first 10,000 integers, we get the following:

```
>>>for i in range(5):
    print("Sum is %d required %10.7f seconds"%sumOfN(10000))
Sum is 50005000 required  0.0018950 seconds
Sum is 50005000 required  0.0018620 seconds
Sum is 50005000 required  0.0019171 seconds
Sum is 50005000 required  0.0019162 seconds
Sum is 50005000 required  0.0019360 seconds
```

We discover that the time is fairly consistent and it takes on average about 0.0019 seconds to execute that code. What if we run the function adding the first 100,000 integers?

```
>>>for i in range(5):
    print("Sum is %d required %10.7f seconds"%sumOfN(100000))
Sum is 5000050000 required  0.0199420 seconds
Sum is 5000050000 required  0.0180972 seconds
Sum is 5000050000 required  0.0194821 seconds
Sum is 5000050000 required  0.0178988 seconds
Sum is 5000050000 required  0.0188949 seconds
```

```
>>>
```

Again, the time required for each run, although longer, is very consistent, averaging about 10 times more seconds. For n equal to 1,000,000 we get:

```
>>>for i in range(5):
    print("Sum is %d required %10.7f seconds"%sumOfN(1000000))
Sum is 500000500000 required  0.1948988 seconds
Sum is 500000500000 required  0.1850290 seconds
Sum is 500000500000 required  0.1809771 seconds
Sum is 500000500000 required  0.1729250 seconds
Sum is 500000500000 required  0.1646299 seconds
>>>
```

In this case, the average again turns out to be about 10 times the previous.

Now consider ActiveCode 3, which shows a different means of solving the summation problem. This function, sumOfN3, takes advantage of a closed equation  $\sum_{i=1}^n i = (n)(n+1)/2$  to compute the sum of the first n integers without iterating.

```
def sumOfN3(n):
    return (n*(n+1))/2

print(sumOfN3(10))
```

代码 17 Summation Without Iteration (active3)

If we do the same benchmark measurement for sumOfN3, using five different values for n (10,000, 100,000, 1,000,000, 10,000,000, and 100,000,000), we get the following results:

```
Sum is 50005000 required 0.00000095 seconds
Sum is 5000050000 required 0.00000191 seconds
Sum is 500000500000 required 0.00000095 seconds
Sum is 50000005000000 required 0.00000095 seconds
Sum is 5000000050000000 required 0.00000119 seconds
```

There are two important things to notice about this output. First, the times recorded above are shorter than any of the previous examples. Second, they are very consistent no matter what the value of n. It appears that sumOfN3 is hardly impacted by the number of integers being added.

But what does this benchmark really tell us? Intuitively, we can see that the iterative solutions seem to be doing more work since some program steps are being repeated. This is likely the reason it is taking longer. Also, the time required for the iterative solution seems to increase as we

increase the value of  $n$ . However, there is a problem. If we ran the same function on a different computer or used a different programming language, we would likely get different results. It could take even longer to perform `sumOfN3` if the computer were older.

We need a better way to characterize these algorithms with respect to execution time. The benchmark technique computes the actual time to execute. It does not really provide us with a useful measurement, because it is dependent on a particular machine, program, time of day, compiler, and programming language. Instead, we would like to have a characterization that is independent of the program or computer being used. This measure would then be useful for judging the algorithm alone and could be used to compare algorithms across implementations.

### 2.2.1. “大 O” 表示法

When trying to characterize an algorithm's efficiency in terms of execution time, independent of any particular program or computer, it is important to quantify the number of operations or steps that the algorithm will require. If each of these steps is considered to be a basic unit of computation, then the execution time for an algorithm can be expressed as the number of steps required to solve the problem. Deciding on an appropriate basic unit of computation can be a complicated problem and will depend on how the algorithm is implemented.

A good basic unit of computation for comparing the summation algorithms shown earlier might be to count the number of assignment statements performed to compute the sum. In the function `sumOfN`, the number of assignment statements is 1 (`theSum=0`) plus the value of  $n$  (the number of times we perform `theSum=theSum+i`). We can denote this by a function, call it  $T$ , where  $T(n)=1+n$ . The parameter  $n$  is often referred to as the “size of the problem,” and we can read this as “ $T(n)$  is the time it takes to solve a problem of size  $n$ , namely  $1+n$  steps.”

In the summation functions given above, it makes sense to use the number of terms in the summation to denote the size of the problem. We can then say that the sum of the first 100,000 integers is a bigger instance of the summation problem than the sum of the first 1,000. Because of this, it might seem reasonable that the time required to solve the larger case would be greater than for the smaller case. Our goal then is to show how the algorithm's execution time changes with respect to the size of the problem.

Computer scientists prefer to take this analysis technique one step further. It turns out that the exact number of operations is not as important as determining the most dominant part of the  $T(n)$  function. In other words, as the problem gets larger, some portion of the  $T(n)$  function tends to overpower the rest. This dominant term is what, in the end, is used for comparison. The order of magnitude function describes the part of  $T(n)$  that increases the fastest as the value of  $n$  increases. Order of magnitude is often called Big-O notation (for “order”) and written as  $O(f(n))$ . It provides a useful approximation to the actual number of steps in the computation. The function  $f(n)$  provides a simple representation of the dominant part of the original  $T(n)$ .

In the above example,  $T(n)=1+n$ . As  $n$  gets large, the constant 1 will become less and less

significant to the final result. If we are looking for an approximation for  $T(n)$ , then we can drop the 1 and simply say that the running time is  $O(n)$ . It is important to note that the 1 is certainly significant for  $T(n)$ . However, as  $n$  gets large, our approximation will be just as accurate without it.

As another example, suppose that for some algorithm, the exact number of steps is  $T(n) = 5n^2 + 27n + 1005$ . When  $n$  is small, say 1 or 2, the constant 1005 seems to be the dominant part of the function. However, as  $n$  gets larger, the  $n^2$  term becomes the most important. In fact, when  $n$  is really large, the other two terms become insignificant in the role that they play in determining the final result. Again, to approximate  $T(n)$  as  $n$  gets large, we can ignore the other terms and focus on  $5n^2$ . In addition, the coefficient 5 becomes insignificant as  $n$  gets large. We would say then that the function  $T(n)$  has an order of magnitude  $f(n) = n^2$ , or simply that it is  $O(n^2)$ .

Although we do not see this in the summation example, sometimes the performance of an algorithm depends on the exact values of the data rather than simply the size of the problem. For these kinds of algorithms we need to characterize their performance in terms of best case, worst case, or average case performance. The worst case performance refers to a particular data set where the algorithm performs especially poorly. Whereas a different data set for the exact same algorithm might have extraordinarily good performance. However, in most cases the algorithm performs somewhere in between these two extremes (average case). It is important for a computer scientist to understand these distinctions so they are not misled by one particular case.

A number of very common order of magnitude functions will come up over and over as you study algorithms. These are shown in Table 1. In order to decide which of these functions is the dominant part of any  $T(n)$  function, we must see how they compare with one another as  $n$  gets large.

Table 1:

$f(n)$	Name
1	Constant
$\log n$	Logarithmic
$n$	Linear
$n \log n$	Log Linear
$n^2$	Quadratic
$n^3$	Cubic
$2^n$	Exponential

表格 11 Common Functions for Big-O

Figure 1 shows graphs of the common functions from Table 1. Notice that when  $n$  is small, the functions are not very well defined with respect to one another. It is hard to tell which is dominant. However, as  $n$  grows, there is a definite relationship and it is easy to see how they compare with one another.

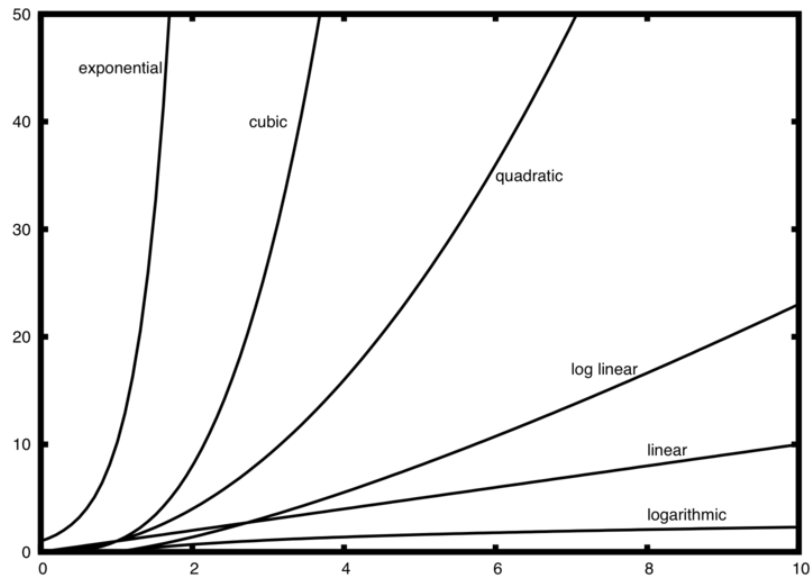


图 13 Plot of Common Big-O Functions

As a final example, suppose that we have the fragment of Python code shown in Listing 2. Although this program does not really do anything, it is instructive to see how we can take actual code and analyze performance.

### Listing 2

```
a=5
b=6
c=10
for i in range(n):
    for j in range(n):
        x = i * i
        y = j * j
        z = i * j
    for k in range(n):
        w = a*k + 45
        v = b*b
d = 33
```

The number of assignment operations is the sum of four terms. The first term is the constant 3, representing the three assignment statements at the start of the fragment. The second term is  $(3n^2)$ , since there are three statements that are performed  $(n^2)$  times due to the nested iteration. The third term is  $(2n)$ , two statements iterated  $n$  times. Finally, the fourth term is the constant 1, representing the final assignment statement. This gives us  $(T(n)=3+3n^2+2n+1=3n^2+2n+4)$ . By looking at the exponents, we can easily see that the  $(n^2)$  term will be dominant and therefore this fragment of code is  $(O(n^2))$ . Note that all of the other terms as well as the coefficient on the dominant term can be ignored as  $n$  grows larger.

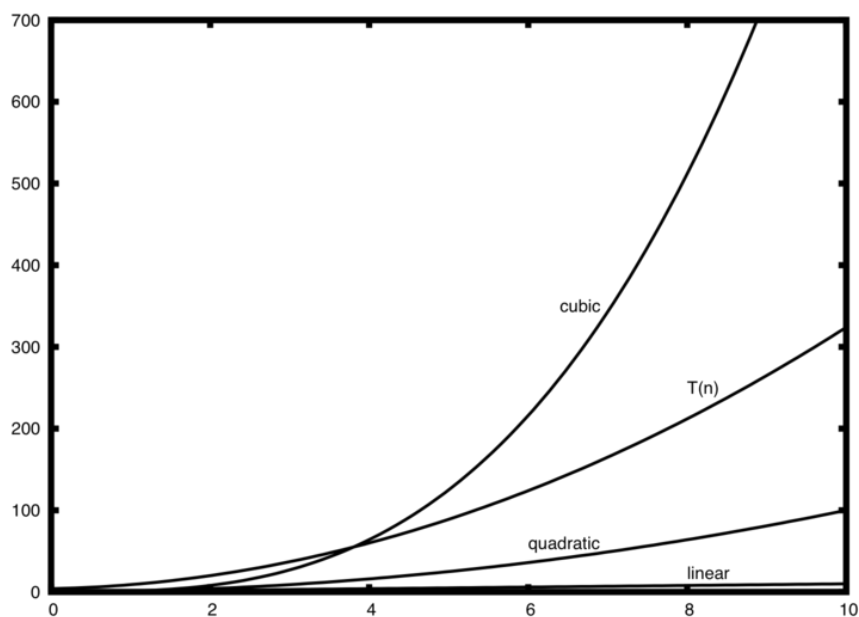


图 14 Comparing  $T(n)$  with Common Big-O Functions

Figure 2 shows a few of the common Big-O functions as they compare with the  $T(n)$  function discussed above. Note that  $T(n)$  is initially larger than the cubic function. However, as  $n$  grows, the cubic function quickly overtakes  $T(n)$ . It is easy to see that  $T(n)$  then follows the quadratic function as  $n$  continues to grow.

#### Self Check

Write two Python functions to find the minimum number in a list. The first function should compare each number to every other number on the list.  $O(n^2)$ . The second function should be linear  $O(n)$ .

### 2.2.2. 例子：“变位词”判断

A good example problem for showing algorithms with different orders of magnitude is the classic anagram detection problem for strings. One string is an anagram of another if the second is simply a rearrangement of the first. For example, 'heart' and 'earth' are anagrams. The strings 'python' and 'typhon' are anagrams as well. For the sake of simplicity, we will assume that the two strings in question are of equal length and that they are made up of symbols from the set of 26 lowercase alphabetic characters. Our goal is to write a boolean function that will take two strings and return whether they are anagrams.

#### 2.2.2.1. 解法 1：检查标记 Checking Off

Our first solution to the anagram problem will check to see that each character in the first string actually occurs in the second. If it is possible to “checkoff” each character, then the two strings must be anagrams. Checking off a character will be accomplished by replacing it with the special



Python value None. However, since strings in Python are immutable, the first step in the process will be to convert the second string to a list. Each character from the first string can be checked against the characters in the list and if found, checked off by replacement. ActiveCode 1 shows this function.

```
def anagramSolution1(s1,s2):
    alist = list(s2)

    pos1 = 0
    stillOK = True

    while pos1 < len(s1) and stillOK:
        pos2 = 0
        found = False
        while pos2 < len(alist) and not found:
            if s1[pos1] == alist[pos2]:
                found = True
            else:
                pos2 = pos2 + 1

        if found:
            alist[pos2] = None
        else:
            stillOK = False

        pos1 = pos1 + 1

    return stillOK

print(anagramSolution1('abcd','dcba'))
```

代码 18 Checking Off (active5)

To analyze this algorithm, we need to note that each of the  $n$  characters in  $s1$  will cause an iteration through up to  $n$  characters in the list from  $s2$ . Each of the  $n$  positions in the list will be visited once to match a character from  $s1$ . The number of visits then becomes the sum of the integers from 1 to  $n$ . We stated earlier that this can be written as

$$\begin{aligned}\sum_{i=1}^n i &= \frac{n(n+1)}{2} \\ &= \frac{1}{2}n^2 + \frac{1}{2}n\end{aligned}$$

As  $n$  gets large, the  $n^2$  term will dominate the  $n$  term and the  $\frac{1}{2}$  can be ignored. Therefore, this solution is  $O(n^2)$ .

### 2.2.2.2. 解法 2：排序比较

Another solution to the anagram problem will make use of the fact that even though `s1` and `s2` are different, they are anagrams only if they consist of exactly the same characters. So, if we begin by sorting each string alphabetically, from a to z, we will end up with the same string if the original two strings are anagrams. ActiveCode 2 shows this solution. Again, in Python we can use the built-in sort method on lists by simply converting each string to a list at the start.

```
def anagramSolution2(s1,s2):
    alist1 = list(s1)
    alist2 = list(s2)

    alist1.sort()
    alist2.sort()

    pos = 0
    matches = True

    while pos < len(s1) and matches:
        if alist1[pos]==alist2[pos]:
            pos = pos + 1
        else:
            matches = False

    return matches

print(anagramSolution2('abcde','edcba'))
```

代码 19 Sort and Compare (active6)

At first glance you may be tempted to think that this algorithm is  $O(n)$ , since there is one simple iteration to compare the  $n$  characters after the sorting process. However, the two calls to the Python sort method are not without their own cost. As we will see in a later chapter, sorting is typically either  $O(n^2)$  or  $O(n\log n)$ , so the sorting operations dominate the iteration. In the end, this algorithm will have the same order of magnitude as that of the sorting process.

### 2.2.2.3. 解法 3：暴力

A brute force technique for solving a problem typically tries to exhaust all possibilities. For the anagram detection problem, we can simply generate a list of all possible strings using the characters from `s1` and then see if `s2` occurs. However, there is a difficulty with this approach. When generating all possible strings from `s1`, there are  $n$  possible first characters,  $n-1$  possible characters for the second position,  $n-2$  for the third, and so on. The total number of candidate strings is  $n * (n-1) * (n-2) * \dots * 3 * 2 * 1$ , which is  $n!$ . Although some of the strings may be

duplicates, the program cannot know this ahead of time and so it will still generate  $n!$  different strings.

It turns out that  $n!$  grows even faster than  $2^n$  as  $n$  gets large. In fact, if  $s1$  were 20 characters long, there would be  $20!=2,432,902,008,176,640,000$  possible candidate strings. If we processed one possibility every second, it would still take us 77,146,816,596 years to go through the entire list. This is probably not going to be a good solution.

#### 2.2.2.4. 解法 4：计数比较

Our final solution to the anagram problem takes advantage of the fact that any two anagrams will have the same number of  $a$ 's, the same number of  $b$ 's, the same number of  $c$ 's, and so on. In order to decide whether two strings are anagrams, we will first count the number of times each character occurs. Since there are 26 possible characters, we can use a list of 26 counters, one for each possible character. Each time we see a particular character, we will increment the counter at that position. In the end, if the two lists of counters are identical, the strings must be anagrams. ActiveCode 3 shows this solution.

```
def anagramSolution4(s1,s2):
    c1 = [0]*26
    c2 = [0]*26

    for i in range(len(s1)):
        pos = ord(s1[i])-ord('a')
        c1[pos] = c1[pos] + 1

    for i in range(len(s2)):
        pos = ord(s2[i])-ord('a')
        c2[pos] = c2[pos] + 1

    j = 0
    stillOK = True
    while j<26 and stillOK:
        if c1[j]==c2[j]:
            j = j + 1
        else:
            stillOK = False

    return stillOK

print(anagramSolution4('apple','pleap'))
```

代码 20 Count and Compare (active7)

Again, the solution has a number of iterations. However, unlike the first solution, none of them are nested. The first two iterations used to count the characters are both based on  $n$ . The third iteration, comparing the two lists of counts, always takes 26 steps since there are 26 possible characters in the strings. Adding it all up gives us  $T(n)=2n+26$  steps. That is  $O(n)$ . We have found a linear order of magnitude algorithm for solving this problem.

Before leaving this example, we need to say something about space requirements. Although the last solution was able to run in linear time, it could only do so by using additional storage to keep the two lists of character counts. In other words, this algorithm sacrificed space in order to gain time.

This is a common occurrence. On many occasions you will need to make decisions between time and space trade-offs. In this case, the amount of extra space is not significant. However, if the underlying alphabet had millions of characters, there would be more concern. As a computer scientist, when given a choice of algorithms, it will be up to you to determine the best use of computing resources given a particular problem.

### Self Check

Q-1: Given the following code fragment, what is its Big-O running time?

```
test = 0
for i in range(n):
    for j in range(n):
        test = test + i * j
```

- a)  $O(n)$
- b)  $O(n^2)$
- c)  $O(\log n)$
- d)  $O(n^3)$

Q-2: Given the following code fragment what is its Big-O running time?

```
test = 0
for i in range(n):
    test = test + 1

for j in range(n):
    test = test - 1
```

- a)  $O(n)$
- b)  $O(n^2)$
- c)  $O(\log n)$
- d)  $O(n^3)$

Q-3: Given the following code fragment what is its Big-O running time?

```
i = n
while i > 0:
```

```
k = 2 + 2  
i = i // 2
```

- a)  $O(n)$
- b)  $O(n^2)$
- c)  $O(\log n)$
- d)  $O(n^3)$

## 2.3. Python 数据结构的性能

Now that you have a general idea of Big-O notation and the differences between the different functions, our goal in this section is to tell you about the Big-O performance for the operations on Python lists and dictionaries. We will then show you some timing experiments that illustrate the costs and benefits of using certain operations on each data structure. It is important for you to understand the efficiency of these Python data structures because they are the building blocks we will use as we implement other data structures in the remainder of the book. In this section we are not going to explain why the performance is what it is. In later chapters you will see some possible implementations of both lists and dictionaries and how the performance depends on the implementation.

### 2.3.1. 列表 List

The designers of Python had many choices to make when they implemented the list data structure. Each of these choices could have an impact on how fast list operations perform. To help them make the right choices they looked at the ways that people would most commonly use the list data structure and they optimized their implementation of a list so that the most common operations were very fast. Of course they also tried to make the less common operations fast, but when a tradeoff had to be made the performance of a less common operation was often sacrificed in favor of the more common operation.

Two common operations are indexing and assigning to an index position. Both of these operations take the same amount of time no matter how large the list becomes. When an operation like this is independent of the size of the list they are  $O(1)$ .

Another very common programming task is to grow a list. There are two ways to create a longer list. You can use the append method or the concatenation operator. The append method is  $O(1)$ . However, the concatenation operator is  $O(k)$  where  $k$  is the size of the list that is being concatenated. This is important for you to know because it can help you make your own programs more efficient by choosing the right tool for the job.

Let's look at four different ways we might generate a list of  $n$  numbers starting with 0. First we'

ll try a for loop and create the list by concatenation, then we ' ll use append rather than concatenation. Next, we' ll try creating the list using list comprehension and finally, and perhaps the most obvious way, using the range function wrapped by a call to the list constructor. Listing 3 shows the code for making our list four different ways.

### Listing 3

```
def test1():
    l = []
    for i in range(1000):
        l = l + [i]

def test2():
    l = []
    for i in range(1000):
        l.append(i)

def test3():
    l = [i for i in range(1000)]

def test4():
    l = list(range(1000))
```

To capture the time it takes for each of our functions to execute we will use Python's timeit module. The timeit module is designed to allow Python developers to make cross-platform timing measurements by running functions in a consistent environment and using timing mechanisms that are as similar as possible across operating systems.

To use timeit you create a Timer object whose parameters are two Python statements. The first parameter is a Python statement that you want to time; the second parameter is a statement that will run once to set up the test. The timeit module will then time how long it takes to execute the statement some number of times. By default timeit will try to run the statement one million times. When its done it returns the time as a floating point value representing the total number of seconds. However, since it executes the statement a million times you can read the result as the number of microseconds to execute the test one time. You can also pass timeit a named parameter called number that allows you to specify how many times the test statement is executed. The following session shows how long it takes to run each of our test functions 1000 times.

```
t1 = Timer("test1()", "from __main__ import test1")
print("concat ",t1.timeit(number=1000), "milliseconds")
t2 = Timer("test2()", "from __main__ import test2")
print("append ",t2.timeit(number=1000), "milliseconds")
t3 = Timer("test3()", "from __main__ import test3")
```

```
print("comprehension ",t3.timeit(number=1000), "milliseconds")
t4 = Timer("test4()", "from __main__ import test4")
print("list range ",t4.timeit(number=1000), "milliseconds")

concat 6.54352807999 milliseconds
append 0.306292057037 milliseconds
comprehension 0.147661924362 milliseconds
list range 0.0655000209808 milliseconds
```

In the experiment above the statement that we are timing is the function call to `test1()`, `test2()`, and so on. The setup statement may look very strange to you, so let's consider it in more detail. You are probably very familiar with the `from`, `import` statement, but this is usually used at the beginning of a Python program file. In this case the statement `from __main__ import test1` imports the function `test1` from the `__main__` namespace into the namespace that `timeit` sets up for the timing experiment. The `timeit` module does this because it wants to run the timing tests in an environment that is uncluttered by any stray variables you may have created, that may interfere with your function's performance in some unforeseen way.

From the experiment above it is clear that the `append` operation at 0.30 milliseconds is much faster than concatenation at 6.54 milliseconds. In the above experiment we also show the times for two additional methods for creating a list; using the list constructor with a call to `range` and a list comprehension. It is interesting to note that the list comprehension is twice as fast as a `for` loop with an `append` operation.

One final observation about this little experiment is that all of the times that you see above include some overhead for actually calling the test function, but we can assume that the function call overhead is identical in all four cases so we still get a meaningful comparison of the operations. So it would not be accurate to say that the concatenation operation takes 6.54 milliseconds but rather the concatenation test function takes 6.54 milliseconds. As an exercise you could test the time it takes to call an empty function and subtract that from the numbers above.

Now that we have seen how performance can be measured concretely you can look at Table 2 to see the Big-O efficiency of all the basic list operations. After thinking carefully about Table 2, you may be wondering about the two different times for `pop`. When `pop` is called on the end of the list it takes  $O(1)$  but when `pop` is called on the first element in the list or anywhere in the middle it is  $O(n)$ . The reason for this lies in how Python chooses to implement lists. When an item is taken from the front of the list, in Python's implementation, all the other elements in the list are shifted one position closer to the beginning. This may seem silly to you now, but if you look at Table 2 you will see that this implementation also allows the index operation to be  $O(1)$ . This is a tradeoff that the Python implementors thought was a good one.

Operation	Big-O Efficiency
<code>index []</code>	$O(1)$

index assignment	O(1)
pop()	O(1)
pop(i)	O(n)
insert(i,item)	O(n)
del operator	O(n)
iteration	O(n)
contains (in)	O(n)
get slice [x:y]	O(k)
del slice	O(n)
set slice	O(n+k)
reverse	O(n)
concatenate	O(k)
sort	O(n log n)
multiply	O(nk)

表格 12 Big-O Efficiency of Python List Operators

As a way of demonstrating this difference in performance let's do another experiment using the `timeit` module. Our goal is to be able to verify the performance of the pop operation on a list of a known size when the program pops from the end of the list, and again when the program pops from the beginning of the list. We will also want to measure this time for lists of different sizes. What we would expect to see is that the time required to pop from the end of the list will stay constant even as the list grows in size, while the time to pop from the beginning of the list will continue to increase as the list grows.

Listing 4 shows one attempt to measure the difference between the two uses of pop. As you can see from this first example, popping from the end takes 0.0003 milliseconds, whereas popping from the beginning takes 4.82 milliseconds. For a list of two million elements this is a factor of 16,000.

There are a couple of things to notice about Listing 4. The first is the statement from `__main__` import x. Although we did not define a function we do want to be able to use the list object x in our test. This approach allows us to time just the single pop statement and get the most accurate measure of the time for that single operation. Because the timer repeats 1000 times it is also important to point out that the list is decreasing in size by 1 each time through the loop. But since the initial list is two million elements in size we only reduce the overall size by 0.05%

#### Listing 4

```
popzero = timeit.Timer("x.pop(0)",
                       "from __main__ import x")
popend = timeit.Timer("x.pop()",
                      "from __main__ import x")
```



```
x = list(range(2000000))
popzero.timeit(number=1000)
4.8213560581207275

x = list(range(2000000))
popend.timeit(number=1000)
0.0003161430358886719
```

While our first test does show that `pop(0)` is indeed slower than `pop()`, it does not validate the claim that `pop(0)` is  $O(n)$  while `pop()` is  $O(1)$ . To validate that claim we need to look at the performance of both calls over a range of list sizes. Listing 5 implements this test.

#### Listing 5

```
popzero = Timer("x.pop(0)",
                "from __main__ import x")
popend = Timer("x.pop()",
               "from __main__ import x")
print("pop(0)    pop()")
for i in range(1000000,100000001,1000000):
    x = list(range(i))
    pt = popend.timeit(number=1000)
    x = list(range(i))
    pz = popzero.timeit(number=1000)
    print("%15.5f, %15.5f" %(pz,pt))
```

Figure 3 shows the results of our experiment. You can see that as the list gets longer and longer the time it takes to `pop(0)` also increases while the time for `pop` stays very flat. This is exactly what we would expect to see for a  $O(n)$  and  $O(1)$  algorithm.

Some sources of error in our little experiment include the fact that there are other processes running on the computer as we measure that may slow down our code, so even though we try to minimize other things happening on the computer there is bound to be some variation in time. That is why the loop runs the test one thousand times in the first place to statistically gather enough information to make the measurement reliable.

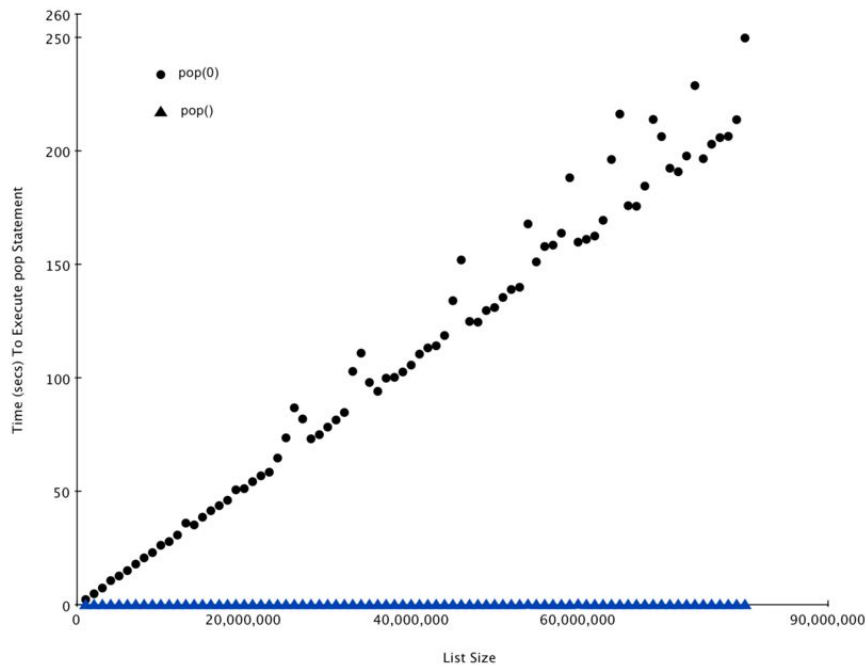


图 15 Comparing the Performance of pop and pop(0)

### 2.3.2. 字典 Dictionary

The second major Python data structure is the dictionary. As you probably recall, dictionaries differ from lists in that you can access items in a dictionary by a key rather than a position. Later in this book you will see that there are many ways to implement a dictionary. The thing that is most important to notice right now is that the get item and set item operations on a dictionary are  $O(1)$ . Another important dictionary operation is the contains operation. Checking to see whether a key is in the dictionary or not is also  $O(1)$ . The efficiency of all dictionary operations is summarized in Table 3. One important side note on dictionary performance is that the efficiencies we provide in the table are for average performance. In some rare cases the contains, get item, and set item operations can degenerate into  $O(n)$  performance but we will get into that in a later chapter when we talk about the different ways that a dictionary could be implemented.

Table 3: Big-O Efficiency of Python Dictionary Operations

operation Big-O Efficiency

copy  $O(n)$

get item  $O(1)$

set item  $O(1)$

delete item  $O(1)$

contains (in)  $O(1)$

iteration  $O(n)$

For our last performance experiment we will compare the performance of the contains operation between lists and dictionaries. In the process we will confirm that the contains operator for lists is  $O(n)$  and the contains operator for dictionaries is  $O(1)$ . The experiment we will use to compare the two is simple. We'll make a list with a range of numbers in it. Then we will pick numbers at

random and check to see if the numbers are in the list. If our performance tables are correct the bigger the list the longer it should take to determine if any one number is contained in the list.

We will repeat the same experiment for a dictionary that contains numbers as the keys. In this experiment we should see that determining whether or not a number is in the dictionary is not only much faster, but the time it takes to check should remain constant even as the dictionary grows larger.

Listing 6 implements this comparison. Notice that we are performing exactly the same operation, `number in container`. The difference is that on line 7 `x` is a list, and on line 9 `x` is a dictionary.

#### Listing 6

```
1 import timeit
2 import random
3
4 for i in range(10000,1000001,20000):
5     t = timeit.Timer("random.randrange(%d) in x"%i,
6                      "from __main__ import random,x")
7     x = list(range(i))
8     lst_time = t.timeit(number=1000)
9     x = {j:None for j in range(i)}
10    d_time = t.timeit(number=1000)
11    print("%d,%10.3f,%10.3f" % (i, lst_time, d_time))
```

Figure 4 summarizes the results of running Listing 6. You can see that the dictionary is consistently faster. For the smallest list size of 10,000 elements a dictionary is 89.4 times faster than a list. For the largest list size of 990,000 elements the dictionary is 11,603 times faster! You can also see that the time it takes for the contains operator on the list grows linearly with the size of the list. This verifies the assertion that the contains operator on a list is  $O(n)$ . It can also be seen that the time for the contains operator on a dictionary is constant even as the dictionary size grows. In fact for a dictionary size of 10,000 the contains operation took 0.004 milliseconds and for the dictionary size of 990,000 it also took 0.004 milliseconds.

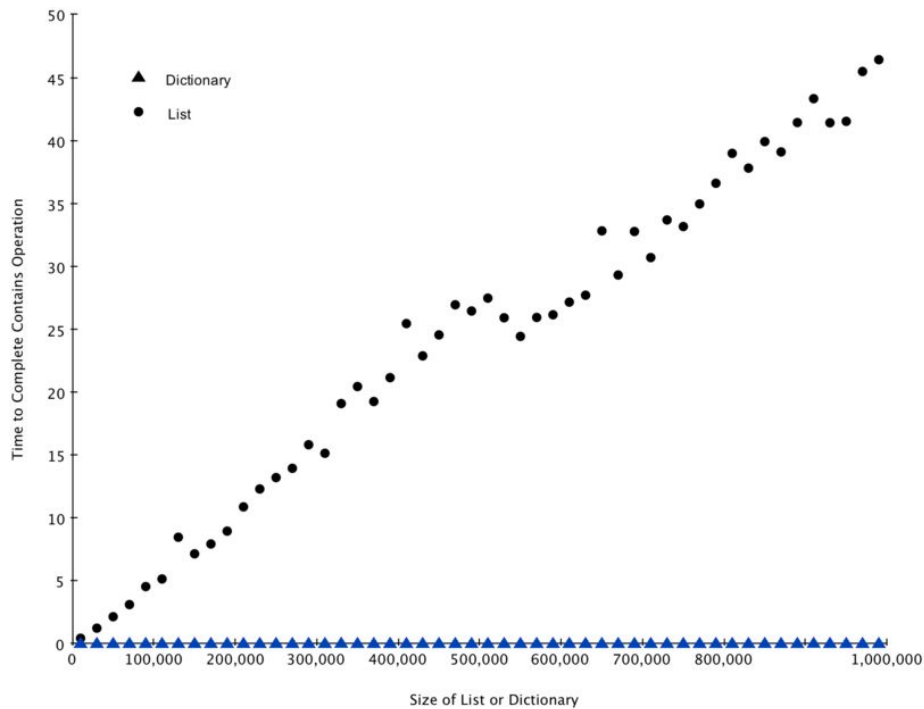


图 16 Comparing the in Operator for Python Lists and Dictionaries

Since Python is an evolving language, there are always changes going on behind the scenes. The latest information on the performance of Python data structures can be found on the Python website. As of this writing the Python wiki has a nice time complexity page that can be found at the Time Complexity Wiki.

### Self Check

Q-4: Which of the list operations shown below is not  $O(1)$ ?

- a) `list.pop(0)`
- b) `list.pop()`
- c) `list.append()`
- d) `list[10]`
- e) all of the above are  $O(1)$

Q-5: Which of the dictionary operations shown below is  $O(1)$ ?

- a) `'x' in mydict`
- b) `del mydict['x']`
- c) `mydict['x'] == 10`
- d) `mydict['x'] = mydict['x'] + 1`
- e) all of the above are  $O(1)$

## 2.4. 小结

- 算法分析是对一个算法进行与具体实现无关的性能分析；
- 大 O 表示法可以依据一个算法对不同规模问题的主干处理过程，对算法的性能进行分

级。

## 2.5. 关键词

average case	Big-O notation	brute force
checking off	exponential	linear
log linear	logarithmic	order of magnitude
quadratic	time complexity	worst case

## 2.6. 问题讨论

1. Give the Big-O performance of the following code fragment:

```
for i in range(n):  
    for j in range(n):  
        k = 2 + 2
```

2. Give the Big-O performance of the following code fragment:

```
for i in range(n):  
    k = 2 + 2
```

3. Give the Big-O performance of the following code fragment:

```
i = n  
while i > 0:  
    k = 2 + 2  
    i = i // 2
```

4. Give the Big-O performance of the following code fragment:

```
for i in range(n):  
    for j in range(n):  
        for k in range(n):  
            k = 2 + 2
```

5. Give the Big-O performance of the following code fragment:

```
i = n  
while i > 0:
```

```
k = 2 + 2  
i = i // 2
```

6. Give the Big-O performance of the following code fragment:

```
for i in range(n):  
    k = 2 + 2  
for j in range(n):  
    k = 2 + 2  
for k in range(n):  
    k = 2 + 2
```

## 2.7. 编程练习

1. Devise an experiment to verify that the list index operator is  $O(1)$
2. Devise an experiment to verify that get item and set item are  $O(1)$  for dictionaries.
3. Devise an experiment that compares the performance of the del operator on lists and dictionaries.
4. Given a list of numbers in random order, write an algorithm that works in  $O(n\log(n))$  to find the kth smallest number in the list.
5. Can you improve the algorithm from the previous problem to be linear? Explain.

## 3. 基本数据结构

### 3.1. 目标

- 了解抽象数据类型：栈 stack、队列 queue、双端队列 deque 和列表 list；
- 能够采用 Python 列表数据类型来实现 stack/queue/deque 等抽象数据类型；
- 了解基本线性数据结构各种具体实现算法的性能；
- 了解前缀、中缀和后缀表达式；
- 采用 stack 对后缀表达式进行求值；
- 采用 stack 将中缀表达式转换为后缀表达式；
- 采用 queue 进行基本的点名报数模拟；
- 能够识别问题属性，选用 stack、queue 或者 deque 中更为合适的数据结构；
- 能够通过节点和节点引用的模式，采用链表来实现抽象数据类型 list；
- 能够比较链表实现与 Python 的 list 实现之间的算法性能。

### 3.2. 什么是线性结构 Linear Structure

We will begin our study of data structures by considering four simple but very powerful concepts. Stacks, queues, dequeues, and lists are examples of data collections whose items are ordered depending on how they are added or removed. Once an item is added, it stays in that position relative to the other elements that came before and came after it. Collections such as these are often referred to as linear data structures.

Linear structures can be thought of as having two ends. Sometimes these ends are referred to as the “left” and the “right” or in some cases the “front” and the “rear.” You could also call them the “top” and the “bottom.” The names given to the ends are not significant. What distinguishes one linear structure from another is the way in which items are added and removed, in particular the location where these additions and removals occur. For example, a structure might allow new items to be added at only one end. Some structures might allow items to be removed from either end.

These variations give rise to some of the most useful data structures in computer science. They appear in many algorithms and can be used to solve a variety of important problems.

### 3.3. 栈 Stack

#### 3.3.1. 什么是栈 stack

A stack (sometimes called a “push-down stack”) is an ordered collection of items where the

addition of new items and the removal of existing items always takes place at the same end. This end is commonly referred to as the “top.” The end opposite the top is known as the “base.”

The base of the stack is significant since items stored in the stack that are closer to the base represent those that have been in the stack the longest. The most recently added item is the one that is in position to be removed first. This ordering principle is sometimes called LIFO, last-in first-out. It provides an ordering based on length of time in the collection. Newer items are near the top, while older items are near the base.

Many examples of stacks occur in everyday situations. Almost any cafeteria has a stack of trays or plates where you take the one at the top, uncovering a new tray or plate for the next customer in line. Imagine a stack of books on a desk (Figure 1). The only book whose cover is visible is the one on top. To access others in the stack, we need to remove the ones that are sitting on top of them. Figure 2 shows another stack. This one contains a number of primitive Python data objects.

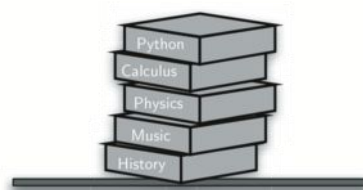


图 17 A Stack of Books

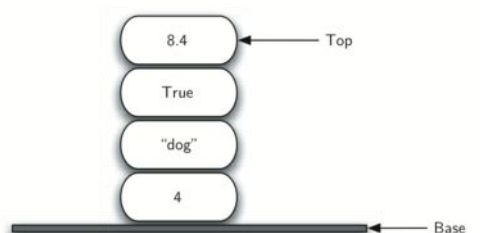


图 18 A Stack of Primitive Python Objects

One of the most useful ideas related to stacks comes from the simple observation of items as they are added and then removed. Assume you start out with a clean desktop. Now place books one at a time on top of each other. You are constructing a stack. Consider what happens when you begin removing books. The order that they are removed is exactly the reverse of the order that they were placed. Stacks are fundamentally important, as they can be used to reverse the order of items. The order of insertion is the reverse of the order of removal. Figure 3 shows the Python data object stack as it was created and then again as items are removed. Note the order of the objects.



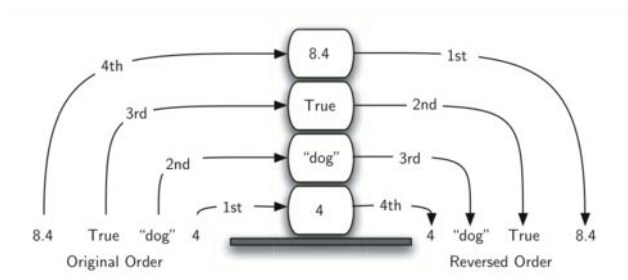


图 19 The Reversal Property of Stacks

Considering this reversal property, you can perhaps think of examples of stacks that occur as you use your computer. For example, every web browser has a Back button. As you navigate from web page to web page, those pages are placed on a stack (actually it is the URLs that are going on the stack). The current page that you are viewing is on the top and the first page you looked at is at the base. If you click on the Back button, you begin to move in reverse order through the pages.

### 3.3.2.抽象数据类型 Stack

The stack abstract data type is defined by the following structure and operations. A stack is structured, as described above, as an ordered collection of items where items are added to and removed from the end called the “top.” Stacks are ordered LIFO. The stack operations are given below.

- `Stack()` creates a new stack that is empty. It needs no parameters and returns an empty stack.
- `push(item)` adds a new item to the top of the stack. It needs the item and returns nothing.
- `pop()` removes the top item from the stack. It needs no parameters and returns the item. The stack is modified.
- `peek()` returns the top item from the stack but does not remove it. It needs no parameters. The stack is not modified.
- `isEmpty()` tests to see whether the stack is empty. It needs no parameters and returns a boolean value.
- `size()` returns the number of items on the stack. It needs no parameters and returns an integer.

For example, if `s` is a stack that has been created and starts out empty, then Table 1 shows the results of a sequence of stack operations. Under stack contents, the top item is listed at the far right.

Stack Operation	Stack Contents	Return Value
<code>s = Stack()</code>	<code>[]</code>	Stack object
<code>s.isEmpty()</code>	<code>[]</code>	True
<code>s.push(4)</code>	<code>[4]</code>	
<code>s.push('dog')</code>	<code>[4, 'dog']</code>	
<code>s.peek()</code>	<code>[4, 'dog']</code>	'dog'

s.push(True)	[4, 'dog', True]	
s.size()	[4, 'dog', True]	3
s.isEmpty()	[4, 'dog', True]	False
s.push(8.4)	[4, 'dog', True, 8.4]	
s.pop()	[4, 'dog', True]	8.4
s.pop()	[4, 'dog']	True
s.size()	[4, 'dog']	2

表格 13 Sample Stack Operations

### 3.3.3. 用 Python 实现 Stack

Now that we have clearly defined the stack as an abstract data type we will turn our attention to using Python to implement the stack. Recall that when we give an abstract data type a physical implementation we refer to the implementation as a data structure.

As we described in Chapter 1, in Python, as in any object-oriented programming language, the implementation of choice for an abstract data type such as a stack is the creation of a new class. The stack operations are implemented as methods. Further, to implement a stack, which is a collection of elements, it makes sense to utilize the power and simplicity of the primitive collections provided by Python. We will use a list.

Recall that the list class in Python provides an ordered collection mechanism and a set of methods. For example, if we have the list [2,5,3,6,7,4], we need only to decide which end of the list will be considered the top of the stack and which will be the base. Once that decision is made, the operations can be implemented using the list methods such as append and pop.

The following stack implementation (ActiveCode 1) assumes that the end of the list will hold the top element of the stack. As the stack grows (as push operations occur), new items will be added on the end of the list. pop operations will manipulate that same end.

```
class Stack:
    def __init__(self):
        self.items = []

    def isEmpty(self):
        return self.items == []

    def push(self, item):
        self.items.append(item)

    def pop(self):
        return self.items.pop()
```

```
def peek(self):
    return self.items[len(self.items)-1]

def size(self):
    return len(self.items)
```

代码 21 Implementing a Stack class using Python lists (stack\_1ac)

Remember that nothing happens when we click the run button other than the definition of the class. We must create a Stack object and then use it. ActiveCode 2 shows the Stack class in action as we perform the sequence of operations from Table 1. Notice that the definition of the Stack class is imported from the pythonds module.

#### Note

The pythonds module contains implementations of all data structures discussed in this book. It is structured according to the sections: basic, trees, and graphs. The module can be downloaded from [pythonworks.org](http://pythonworks.org).

```
from pythonds.basic.stack import Stack

s=Stack()

print(s.isEmpty())
s.push(4)
s.push('dog')
print(s.peek())
s.push(True)
print(s.size())
print(s.isEmpty())
s.push(8.4)
print(s.pop())
print(s.pop())
print(s.size())
```

代码 22 (stack\_ex\_1)

It is important to note that we could have chosen to implement the stack using a list where the top is at the beginning instead of at the end. In this case, the previous pop and append methods would no longer work and we would have to index position 0 (the first item in the list) explicitly using pop and insert. The implementation is shown in CodeLens 1.

```
1 class Stack:
2     def __init__(self):
```

```

3         self.items = []
4
5     def isEmpty(self):
6         return self.items == []
7
8     def push(self, item):
9         self.items.insert(0,item)
10
11    def pop(self):
12        return self.items.pop(0)
13
14    def peek(self):
15        return self.items[0]
16
17    def size(self):
18        return len(self.items)
19
20 s = Stack()
21 s.push('hello')
22 s.push('true')
23 print(s.pop())

```

代码 23 Alternative Implementation of the Stack class (stack\_cl\_1)

This ability to change the physical implementation of an abstract data type while maintaining the logical characteristics is an example of abstraction at work. However, even though the stack will work either way, if we consider the performance of the two implementations, there is definitely a difference. Recall that the append and pop() operations were both  $O(1)$ . This means that the first implementation will perform push and pop in constant time no matter how many items are on the stack. The performance of the second implementation suffers in that the insert(0) and pop(0) operations will both require  $O(n)$  for a stack of size  $n$ . Clearly, even though the implementations are logically equivalent, they would have very different timings when performing benchmark testing.

### Self Check

Q-10: Given the following sequence of stack operations, what is the top item on the stack when the sequence is complete?

```

m = Stack()
m.push('x')
m.push('y')
m.pop()
m.push('z')
m.peek()

```

a) 'x'

- b) 'y'
- c) 'z'
- d) The stack is empty

Q-11: Given the following sequence of stack operations, what is the top item on the stack when the sequence is complete?

```
m = Stack()
m.push('x')
m.push('y')
m.push('z')
while not m.isEmpty():
    m.pop()
    m.pop()
```

- a) 'x'
- b) the stack is empty
- c) an error will occur
- d) 'z'

Write a function `revstring(mystr)` that uses a stack to reverse the characters in a string.

### 3.3.4. 简单括号匹配

We now turn our attention to using stacks to solve real computer science problems. You have no doubt written arithmetic expressions such as

```
(5+6)* (7+8)/(4+3)
```

where parentheses are used to order the performance of operations. You may also have some experience programming in a language such as Lisp with constructs like

```
(defun square(n)
  (* n n))
```

This defines a function called `square` that will return the square of its argument `n`. Lisp is notorious for using lots and lots of parentheses.

In both of these examples, parentheses must appear in a balanced fashion. Balanced parentheses means that each opening symbol has a corresponding closing symbol and the pairs of parentheses are properly nested. Consider the following correctly balanced strings of parentheses:

((()()()))

((((( )))

((()((()()())))

Compare those with the following, which are not balanced:

((((( )))

( )))

((()()())

The ability to differentiate between parentheses that are correctly balanced and those that are unbalanced is an important part of recognizing many programming language structures.

The challenge then is to write an algorithm that will read a string of parentheses from left to right and decide whether the symbols are balanced. To solve this problem we need to make an important observation. As you process symbols from left to right, the most recent opening parenthesis must match the next closing symbol (see Figure 4). Also, the first opening symbol processed may have to wait until the very last symbol for its match. Closing symbols match opening symbols in the reverse order of their appearance; they match from the inside out. This is a clue that stacks can be used to solve the problem.

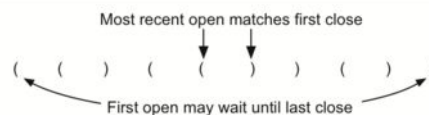


图 20 Matching Parentheses

Once you agree that a stack is the appropriate data structure for keeping the parentheses, the statement of the algorithm is straightforward. Starting with an empty stack, process the parenthesis strings from left to right. If a symbol is an opening parenthesis, push it on the stack as a signal that a corresponding closing symbol needs to appear later. If, on the other hand, a symbol is a closing parenthesis, pop the stack. As long as it is possible to pop the stack to match every closing symbol, the parentheses remain balanced. If at any time there is no opening symbol on the stack to match a closing symbol, the string is not balanced properly. At the end of the string, when all symbols have been processed, the stack should be empty. The Python code to implement this algorithm is shown in ActiveCode 1.

```
from pythonds.basic.stack import Stack

def parChecker(symbolString):
    s = Stack()
    balanced = True
```

```

index = 0
while index < len(symbolString) and balanced:
    symbol = symbolString[index]
    if symbol == "(":
        s.push(symbol)
    else:
        if s.isEmpty():
            balanced = False
        else:
            s.pop()

    index = index + 1

if balanced and s.isEmpty():
    return True
else:
    return False

print(parChecker('((( )))'))
print(parChecker('(() )'))

```

代码 24 Solving the Balanced Parentheses Problem (parcheck1)

This function, `parChecker`, assumes that a `Stack` class is available and returns a boolean result as to whether the string of parentheses is balanced. Note that the boolean variable `balanced` is initialized to `True` as there is no reason to assume otherwise at the start. If the current symbol is `(`, then it is pushed on the stack (lines 9 – 10). Note also in line 15 that `pop` simply removes a symbol from the stack. The returned value is not used since we know it must be an opening symbol seen earlier. At the end (lines 19 – 22), as long as the expression is balanced and the stack has been completely cleaned off, the string represents a correctly balanced sequence of parentheses.

### 3.3.5. 匹配符号（通用情况）

The balanced parentheses problem shown above is a specific case of a more general situation that arises in many programming languages. The general problem of balancing and nesting different kinds of opening and closing symbols occurs frequently. For example, in Python square brackets, `[` and `]`, are used for lists; curly braces, `{` and `}`, are used for dictionaries; and parentheses, `(` and `)`, are used for tuples and arithmetic expressions. It is possible to mix symbols as long as each maintains its own open and close relationship. Strings of symbols such as

```

{ { ( [ ] [ ] ) } ( ) }

[ [ { { ( ( ) ) } } ] ]

```

```
[ ] [ ] [ ] ( ) { }
```

are properly balanced in that not only does each opening symbol have a corresponding closing symbol, but the types of symbols match as well.

Compare those with the following strings that are not balanced:

```
( [ ) ]
```

```
(( ( ) ] ) )
```

```
[ { ( ) ]
```

The simple parentheses checker from the previous section can easily be extended to handle these new types of symbols. Recall that each opening symbol is simply pushed on the stack to wait for the matching closing symbol to appear later in the sequence. When a closing symbol does appear, the only difference is that we must check to be sure that it correctly matches the type of the opening symbol on top of the stack. If the two symbols do not match, the string is not balanced. Once again, if the entire string is processed and nothing is left on the stack, the string is correctly balanced.

The Python program to implement this is shown in ActiveCode 1. The only change appears in line 16 where we call a helper function, `matches`, to assist with symbol-matching. Each symbol that is removed from the stack must be checked to see that it matches the current closing symbol. If a mismatch occurs, the boolean variable `balanced` is set to `False`.

```
from pythonds.basic.stack import Stack

def parChecker(symbolString):
    s = Stack()
    balanced = True
    index = 0
    while index < len(symbolString) and balanced:
        symbol = symbolString[index]
        if symbol in "([{":
            s.push(symbol)
        else:
            if s.isEmpty():
                balanced = False
            else:
                top = s.pop()
                if not matches(top, symbol):
                    balanced = False
        index = index + 1
```



```

    if balanced and s.isEmpty():
        return True
    else:
        return False

def matches(open,close):
    opens = "([{"
    closers = ")]}"
    return opens.index(open) == closers.index(close)

print(parChecker('{{([][])}()})'))
print(parChecker('[{()}]'))

```

代码 25 Solving the General Balanced Symbol Problem (parcheck2)

These two examples show that stacks are very important data structures for the processing of language constructs in computer science. Almost any notation you can think of has some type of nested symbol that must be matched in a balanced order. There are a number of other important uses for stacks in computer science. We will continue to explore them in the next sections.

### 3.3.6. 十进制数转换为二进制

In your study of computer science, you have probably been exposed in one way or another to the idea of a binary number. Binary representation is important in computer science since all values stored within a computer exist as a string of binary digits, a string of 0s and 1s. Without the ability to convert back and forth between common representations and binary numbers, we would need to interact with computers in very awkward ways.

Integer values are common data items. They are used in computer programs and computation all the time. We learn about them in math class and of course represent them using the decimal number system, or base 10. The decimal number  $233_{10}$  and its corresponding binary equivalent  $11101001_2$  are interpreted respectively as

$$2 \times 10^2 + 3 \times 10^1 + 3 \times 10^0$$

and

$$1 \times 2^7 + 1 \times 2^6 + 1 \times 2^5 + 0 \times 2^4 + 1 \times 2^3 + 0 \times 2^2 + 0 \times 2^1 + 1 \times 2^0$$

But how can we easily convert integer values into binary numbers? The answer is an algorithm called “Divide by 2” that uses a stack to keep track of the digits for the binary result.

The Divide by 2 algorithm assumes that we start with an integer greater than 0. A simple iteration then continually divides the decimal number by 2 and keeps track of the remainder. The first division by 2 gives information as to whether the value is even or odd. An even value will have a remainder of 0. It will have the digit 0 in the ones place. An odd value will have a remainder of 1 and will have the digit 1 in the ones place. We think about building our binary number as a sequence of digits; the first remainder we compute will actually be the last digit in the sequence. As shown in Figure 5, we again see the reversal property that signals that a stack is likely to be the appropriate data structure for solving the problem.

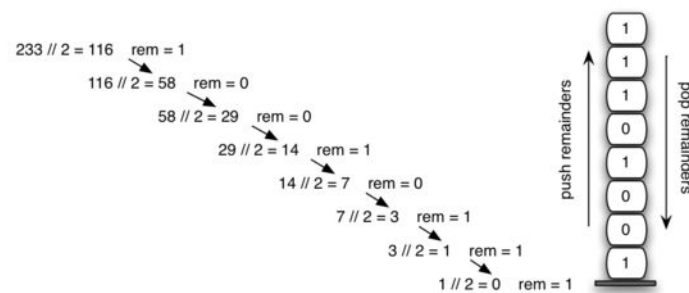


图 21 Decimal-to-Binary Conversion

The Python code in ActiveCode 1 implements the Divide by 2 algorithm. The function `divideBy2` takes an argument that is a decimal number and repeatedly divides it by 2. Line 7 uses the built-in modulo operator, `%`, to extract the remainder and line 8 then pushes it on the stack. After the division process reaches 0, a binary string is constructed in lines 11-13. Line 11 creates an empty string. The binary digits are popped from the stack one at a time and appended to the right-hand end of the string. The binary string is then returned.

```
from pythonds.basic.stack import Stack

def divideBy2(decNumber):
    remstack = Stack()

    while decNumber > 0:
        rem = decNumber % 2
        remstack.push(rem)
        decNumber = decNumber // 2

    binString = ""
    while not remstack.isEmpty():
        binString = binString + str(remstack.pop())

    return binString

print(divideBy2(42))
```

## 代码 26 Converting from Decimal to Binary (divby2)

The algorithm for binary conversion can easily be extended to perform the conversion for any base. In computer science it is common to use a number of different encodings. The most common of these are binary, octal (base 8), and hexadecimal (base 16).

The decimal number 233 and its corresponding octal and hexadecimal equivalents  $351_8$  and  $E9_{16}$  are interpreted as

$$3 \times 8^2 + 5 \times 8^1 + 1 \times 8^0$$

and

$$14 \times 16^1 + 9 \times 16^0$$

The function `divideBy2` can be modified to accept not only a decimal value but also a base for the intended conversion. The “Divide by 2” idea is simply replaced with a more general “Divide by base.” A new function called `baseConverter`, shown in ActiveCode 2, takes a decimal number and any base between 2 and 16 as parameters. The remainders are still pushed onto the stack until the value being converted becomes 0. The same left-to-right string construction technique can be used with one slight change. Base 2 through base 10 numbers need a maximum of 10 digits, so the typical digit characters 0, 1, 2, 3, 4, 5, 6, 7, 8, and 9 work fine. The problem comes when we go beyond base 10. We can no longer simply use the remainders, as they are themselves represented as two-digit decimal numbers. Instead we need to create a set of digits that can be used to represent those remainders beyond 9.

```
from pythonds.basic.stack import Stack

def baseConverter(decNumber, base):
    digits = "0123456789ABCDEF"

    remstack = Stack()

    while decNumber > 0:
        rem = decNumber % base
        remstack.push(rem)
        decNumber = decNumber // base

    newString = ""
    while not remstack.isEmpty():
        newString = newString + digits[remstack.pop()]

    return newString
```

```
print(baseConverter(25,2))  
print(baseConverter(25,16))
```

代码 27 Converting from Decimal to any Base (baseconvert)

A solution to this problem is to extend the digit set to include some alphabet characters. For example, hexadecimal uses the ten decimal digits along with the first six alphabet characters for the 16 digits. To implement this, a digit string is created (line 4 in Listing 6) that stores the digits in their corresponding positions. 0 is at position 0, 1 is at position 1, A is at position 10, B is at position 11, and so on. When a remainder is removed from the stack, it can be used to index into the digit string and the correct resulting digit can be appended to the answer. For example, if the remainder 13 is removed from the stack, the digit D is appended to the resulting string.

#### Self Check

Q-6: What is value of 25 expressed as an octal number

Q-7: What is value of 256 expressed as a hexadecimal number

Q-8: What is value of 26 expressed in base 26

### 3.3.7. 中綴、前綴和后綴表达式

When you write an arithmetic expression such as  $B * C$ , the form of the expression provides you with information so that you can interpret it correctly. In this case we know that the variable  $B$  is being multiplied by the variable  $C$  since the multiplication operator  $*$  appears between them in the expression. This type of notation is referred to as infix since the operator is in between the two operands that it is working on.

Consider another infix example,  $A + B * C$ . The operators  $+$  and  $*$  still appear between the operands, but there is a problem. Which operands do they work on? Does the  $+$  work on  $A$  and  $B$  or does the  $*$  take  $B$  and  $C$ ? The expression seems ambiguous.

In fact, you have been reading and writing these types of expressions for a long time and they do not cause you any problem. The reason for this is that you know something about the operators  $+$  and  $*$ . Each operator has a precedence level. Operators of higher precedence are used before operators of lower precedence. The only thing that can change that order is the presence of parentheses. The precedence order for arithmetic operators places multiplication and division above addition and subtraction. If two operators of equal precedence appear, then a left-to-right ordering or associativity is used.

Let's interpret the troublesome expression  $A + B * C$  using operator precedence.  $B$  and  $C$  are multiplied first, and  $A$  is then added to that result.  $(A + B) * C$  would force the addition of  $A$  and  $B$  to be done first before the multiplication. In expression  $A + B + C$ , by precedence (via

associativity), the leftmost + would be done first.

Although all this may be obvious to you, remember that computers need to know exactly what operators to perform and in what order. One way to write an expression that guarantees there will be no confusion with respect to the order of operations is to create what is called a fully parenthesized expression. This type of expression uses one pair of parentheses for each operator. The parentheses dictate the order of operations; there is no ambiguity. There is also no need to remember any precedence rules.

The expression  $A + B * C + D$  can be rewritten as  $((A + (B * C)) + D)$  to show that the multiplication happens first, followed by the leftmost addition.  $A + B + C + D$  can be written as  $((A + B) + C) + D$  since the addition operations associate from left to right.

There are two other very important expression formats that may not seem obvious to you at first. Consider the infix expression  $A + B$ . What would happen if we moved the operator before the two operands? The resulting expression would be  $+ A B$ . Likewise, we could move the operator to the end. We would get  $A B +$ . These look a bit strange.

These changes to the position of the operator with respect to the operands create two new expression formats, prefix and postfix. Prefix expression notation requires that all operators precede the two operands that they work on. Postfix, on the other hand, requires that its operators come after the corresponding operands. A few more examples should help to make this a bit clearer (see Table 2).

$A + B * C$  would be written as  $+ A * B C$  in prefix. The multiplication operator comes immediately before the operands  $B$  and  $C$ , denoting that  $*$  has precedence over  $+$ . The addition operator then appears before the  $A$  and the result of the multiplication.

In postfix, the expression would be  $A B C * +$ . Again, the order of operations is preserved since the  $*$  appears immediately after the  $B$  and the  $C$ , denoting that  $*$  has precedence, with  $+$  coming after. Although the operators moved and now appear either before or after their respective operands, the order of the operands stayed exactly the same relative to one another.

Infix Expression	Prefix Expression	Postfix Expression
$A + B$	$+ A B$	$A B +$
$A + B * C$	$+ A * B C$	$A B C * +$

表格 14 Examples of Infix, Prefix, and Postfix

Now consider the infix expression  $(A + B) * C$ . Recall that in this case, infix requires the parentheses to force the performance of the addition before the multiplication. However, when  $A + B$  was written in prefix, the addition operator was simply moved before the operands,  $+ A B$ . The result of this operation becomes the first operand for the multiplication. The multiplication operator is moved in front of the entire expression, giving us  $* + A B C$ . Likewise, in postfix  $A B +$  forces the addition to happen first. The multiplication can be done to that result and the

remaining operand C. The proper postfix expression is then  $A B + C *$ .

Consider these three expressions again (see Table 3). Something very important has happened. Where did the parentheses go? Why don't we need them in prefix and postfix? The answer is that the operators are no longer ambiguous with respect to the operands that they work on. Only infix notation requires the additional symbols. The order of operations within prefix and postfix expressions is completely determined by the position of the operator and nothing else. In many ways, this makes infix the least desirable notation to use.

Infix Expression	Prefix Expression	Postfix Expression
$(A + B) * C$	$* + A B C$	$A B + C *$

表格 15 An Expression with Parentheses

Table 4 shows some additional examples of infix expressions and the equivalent prefix and postfix expressions. Be sure that you understand how they are equivalent in terms of the order of the operations being performed.

Infix Expression	Prefix Expression	Postfix Expression
$A + B * C + D$	$+ + A * B C D$	$A B C * + D +$
$(A + B) * (C + D)$	$* + A B + C D$	$A B + C D + *$
$A * B + C * D$	$+ * A B * C D$	$A B * C D * +$
$A + B + C + D$	$+ + + A B C D$	$A B + C + D +$

表格 16 Additional Examples of Infix, Prefix, and Postfix

### 3.3.7.1. 中缀表达式转换为前缀和后缀形式

So far, we have used ad hoc methods to convert between infix expressions and the equivalent prefix and postfix expression notations. As you might expect, there are algorithmic ways to perform the conversion that allow any expression of any complexity to be correctly transformed.

The first technique that we will consider uses the notion of a fully parenthesized expression that was discussed earlier. Recall that  $A + B * C$  can be written as  $(A + (B * C))$  to show explicitly that the multiplication has precedence over the addition. On closer observation, however, you can see that each parenthesis pair also denotes the beginning and the end of an operand pair with the corresponding operator in the middle.

Look at the right parenthesis in the subexpression  $(B * C)$  above. If we were to move the multiplication symbol to that position and remove the matching left parenthesis, giving us  $B C *$ , we would in effect have converted the subexpression to postfix notation. If the addition operator were also moved to its corresponding right parenthesis position and the matching left parenthesis were removed, the complete postfix expression would result (see Figure 6).



图 22 Moving Operators to the Right for Postfix Notation

If we do the same thing but instead of moving the symbol to the position of the right parenthesis, we move it to the left, we get prefix notation (see Figure 7). The position of the parenthesis pair is actually a clue to the final position of the enclosed operator.

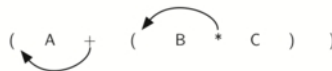


图 23 Moving Operators to the Left for Prefix Notation

So in order to convert an expression, no matter how complex, to either prefix or postfix notation, fully parenthesize the expression using the order of operations. Then move the enclosed operator to the position of either the left or the right parenthesis depending on whether you want prefix or postfix notation.

Here is a more complex expression:  $(A + B) * C - (D - E) * (F + G)$ . Figure 8 shows the conversion to postfix and prefix notations.

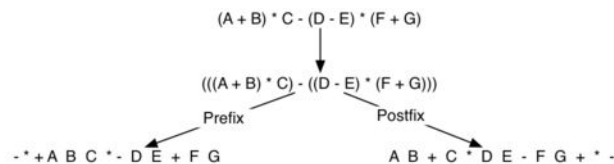


图 24 Converting a Complex Expression to Prefix and Postfix Notations

### 3.3.7.2. 通用的中缀转后缀算法

We need to develop an algorithm to convert any infix expression to a postfix expression. To do this we will look closer at the conversion process.

Consider once again the expression  $A + B * C$ . As shown above,  $A B C * +$  is the postfix equivalent. We have already noted that the operands  $A$ ,  $B$ , and  $C$  stay in their relative positions. It is only the operators that change position. Let's look again at the operators in the infix expression. The first operator that appears from left to right is  $+$ . However, in the postfix expression,  $+$  is at the end since the next operator,  $*$ , has precedence over addition. The order of the operators in the original expression is reversed in the resulting postfix expression.

As we process the expression, the operators have to be saved somewhere since their corresponding right operands are not seen yet. Also, the order of these saved operators may need to be reversed due to their precedence. This is the case with the addition and the multiplication in this example.

Since the addition operator comes before the multiplication operator and has lower precedence, it needs to appear after the multiplication operator is used. Because of this reversal of order, it makes sense to consider using a stack to keep the operators until they are needed.

What about  $(A + B) * C$ ? Recall that  $A B + C *$  is the postfix equivalent. Again, processing this infix expression from left to right, we see  $+$  first. In this case, when we see  $*$ ,  $+$  has already been placed in the result expression because it has precedence over  $*$  by virtue of the parentheses. We can now start to see how the conversion algorithm will work. When we see a left parenthesis, we will save it to denote that another operator of high precedence will be coming. That operator will need to wait until the corresponding right parenthesis appears to denote its position (recall the fully parenthesized technique). When that right parenthesis does appear, the operator can be popped from the stack.

As we scan the infix expression from left to right, we will use a stack to keep the operators. This will provide the reversal that we noted in the first example. The top of the stack will always be the most recently saved operator. Whenever we read a new operator, we will need to consider how that operator compares in precedence with the operators, if any, already on the stack.

Assume the infix expression is a string of tokens delimited by spaces. The operator tokens are  $*$ ,  $/$ ,  $+$ , and  $-$ , along with the left and right parentheses,  $($  and  $)$ . The operand tokens are the single-character identifiers  $A$ ,  $B$ ,  $C$ , and so on. The following steps will produce a string of tokens in postfix order.

1. Create an empty stack called opstack for keeping operators. Create an empty list for output.
2. Convert the input infix string to a list by using the string method split.
3. Scan the token list from left to right.
  - a) If the token is an operand, append it to the end of the output list.
  - b) If the token is a left parenthesis, push it on the opstack.
  - c) If the token is a right parenthesis, pop the opstack until the corresponding left parenthesis is removed. Append each operator to the end of the output list.
  - d) If the token is an operator,  $*$ ,  $/$ ,  $+$ , or  $-$ , push it on the opstack. However, first remove any operators already on the opstack that have higher or equal precedence and append them to the output list.
4. When the input expression has been completely processed, check the opstack. Any operators still on the stack can be removed and appended to the end of the output list.

Figure 9 shows the conversion algorithm working on the expression  $A * B + C * D$ . Note that the first  $*$  operator is removed upon seeing the  $+$  operator. Also,  $+$  stays on the stack when the second  $*$  occurs, since multiplication has precedence over addition. At the end of the infix expression the stack is popped twice, removing both operators and placing  $+$  as the last operator in the postfix expression.



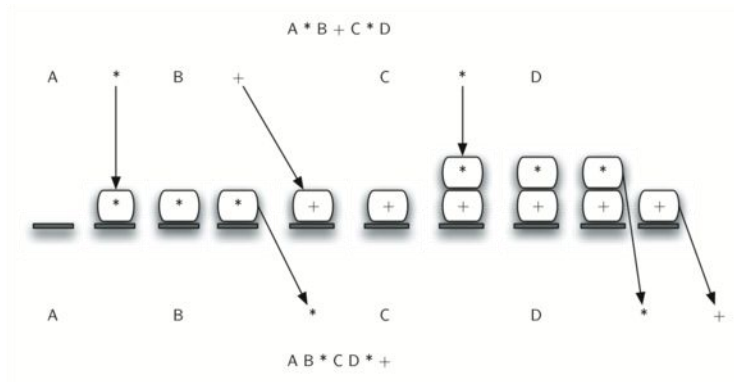


图 25 Converting  $A * B + C * D$  to Postfix Notation

In order to code the algorithm in Python, we will use a dictionary called `prec` to hold the precedence values for the operators. This dictionary will map each operator to an integer that can be compared against the precedence levels of other operators (we have arbitrarily used the integers 3, 2, and 1). The left parenthesis will receive the lowest value possible. This way any operator that is compared against it will have higher precedence and will be placed on top of it. Line 15 defines the operands to be any upper-case character or digit. The complete conversion function is shown in ActiveCode 1.

```
from pythonds.basic.stack import Stack

def infixToPostfix(infixexpr):
    prec = {}
    prec["*"] = 3
    prec["/"] = 3
    prec["+"] = 2
    prec["-"] = 2
    prec["("] = 1
    opStack = Stack()
    postfixList = []
    tokenList = infixexpr.split()

    for token in tokenList:
        if token in "ABCDEFGHIJKLMNOPQRSTUVWXYZ" or token in "0123456789":
            postfixList.append(token)
        elif token == '(':
            opStack.push(token)
        elif token == ')':
            topToken = opStack.pop()
            while topToken != '(':
                postfixList.append(topToken)
                topToken = opStack.pop()
            topToken = opStack.pop()
        else:
            while opStack and prec[opStack[-1]] >= prec[token]:
                postfixList.append(opStack.pop())
            opStack.push(token)
```

```

        while (not opStack.isEmpty()) and \
            (prec[opStack.peek()] >= prec[token]):
            postfixList.append(opStack.pop())
        opStack.push(token)

    while not opStack.isEmpty():
        postfixList.append(opStack.pop())
    return " ".join(postfixList)

print(infixToPostfix("A * B + C * D"))
print(infixToPostfix("( A + B ) * C - ( D - E ) * ( F + G )"))

```

代码 28 Converting Infix Expressions to Postfix Expressions (intopost)

A few more examples of execution in the Python shell are shown below.

```

>>> infixtopostfix("( A + B ) * ( C + D )")
'A B + C D + *'
>>> infixtopostfix("( A + B ) * C")
'A B + C *'
>>> infixtopostfix("A + B * C")
'A B C * +'
>>>

```

### 3.3.7.3. 后缀表达式求值

As a final stack example, we will consider the evaluation of an expression that is already in postfix notation. In this case, a stack is again the data structure of choice. However, as you scan the postfix expression, it is the operands that must wait, not the operators as in the conversion algorithm above. Another way to think about the solution is that whenever an operator is seen on the input, the two most recent operands will be used in the evaluation.

To see this in more detail, consider the postfix expression `4 5 6 * +`. As you scan the expression from left to right, you first encounter the operands 4 and 5. At this point, you are still unsure what to do with them until you see the next symbol. Placing each on the stack ensures that they are available if an operator comes next.

In this case, the next symbol is another operand. So, as before, push it and check the next symbol. Now we see an operator, `*`. This means that the two most recent operands need to be used in a multiplication operation. By popping the stack twice, we can get the proper operands and then perform the multiplication (in this case getting the result 30).

We can now handle this result by placing it back on the stack so that it can be used as an operand for the later operators in the expression. When the final operator is processed, there will be only one value left on the stack. Pop and return it as the result of the expression. Figure 10 shows the stack contents as this entire example expression is being processed.

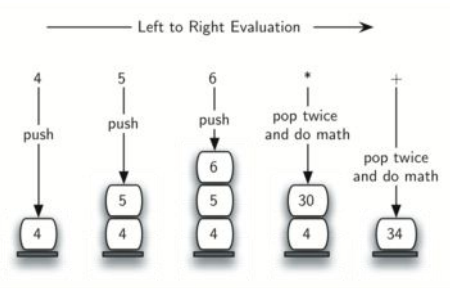


图 26 Stack Contents During Evaluation

Figure 11 shows a slightly more complex example,  $7\ 8\ +\ 3\ 2\ +\ /\$ . There are two things to note in this example. First, the stack size grows, shrinks, and then grows again as the subexpressions are evaluated. Second, the division operation needs to be handled carefully. Recall that the operands in the postfix expression are in their original order since postfix changes only the placement of operators. When the operands for the division are popped from the stack, they are reversed. Since division is not a commutative operator, in other words  $15/5$  is not the same as  $5/15$ , we must be sure that the order of the operands is not switched.

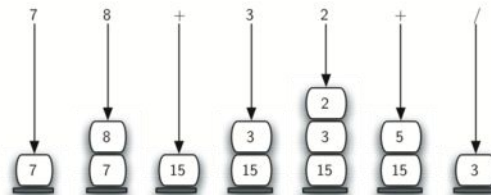


图 27 A More Complex Example of Evaluation

Assume the postfix expression is a string of tokens delimited by spaces. The operators are  $*$ ,  $/$ ,  $+$ , and  $-$  and the operands are assumed to be single-digit integer values. The output will be an integer result.

1. Create an empty stack called operandStack.
2. Convert the string to a list by using the string method split.
3. Scan the token list from left to right.
  - a) If the token is an operand, convert it from a string to an integer and push the value onto the operandStack.
  - b) If the token is an operator,  $*$ ,  $/$ ,  $+$ , or  $-$ , it will need two operands. Pop the operandStack twice. The first pop is the second operand and the second pop is the first operand. Perform the arithmetic operation. Push the result back on the operandStack.
4. When the input expression has been completely processed, the result is on the stack. Pop the operandStack and return the value.

The complete function for the evaluation of postfix expressions is shown in ActiveCode 2. To assist with the arithmetic, a helper function doMath is defined that will take two operands and an operator and then perform the proper arithmetic operation.

```
from pythonds.basic.stack import Stack

def postfixEval(postfixExpr):
    operandStack = Stack()
    tokenList = postfixExpr.split()

    for token in tokenList:
        if token in "0123456789":
            operandStack.push(int(token))
        else:
            operand2 = operandStack.pop()
            operand1 = operandStack.pop()
            result = doMath(token, operand1, operand2)
            operandStack.push(result)
    return operandStack.pop()

def doMath(op, op1, op2):
    if op == "*":
        return op1 * op2
    elif op == "/":
        return op1 / op2
    elif op == "+":
        return op1 + op2
    else:
        return op1 - op2

print(postfixEval('7 8 + 3 2 + /'))
```

代码 29 Postfix Evaluation (postfixeval)

It is important to note that in both the postfix conversion and the postfix evaluation programs we assumed that there were no errors in the input expression. Using these programs as a starting point, you can easily see how error detection and reporting can be included. We leave this as an exercise at the end of the chapter.

### Self Check

Q-12: Without using the activecode infixToPostfix function, convert the following expression to postfix  $10 + 3 * 5 / (16 - 4)$

Q-13:  $17 10 + 3 * 9 / ==$

Q-14: Modify the infixToPostfix function so that it can convert the following expression:  $5 * 3 ^ (4 - 2)$  Paste the answer here:

## 3.4. 队列 queue

### 3.4.1. 什么是队列

A queue is an ordered collection of items where the addition of new items happens at one end, called the “rear,” and the removal of existing items occurs at the other end, commonly called the “front.” As an element enters the queue it starts at the rear and makes its way toward the front, waiting until that time when it is the next element to be removed.

The most recently added item in the queue must wait at the end of the collection. The item that has been in the collection the longest is at the front. This ordering principle is sometimes called FIFO, first-in first-out. It is also known as “first-come first-served.”

The simplest example of a queue is the typical line that we all participate in from time to time. We wait in a line for a movie, we wait in the check-out line at a grocery store, and we wait in the cafeteria line (so that we can pop the tray stack). Well-behaved lines, or queues, are very restrictive in that they have only one way in and only one way out. There is no jumping in the middle and no leaving before you have waited the necessary amount of time to get to the front. Figure 1 shows a simple queue of Python data objects.



图 28 A Queue of Python Data Objects

Computer science also has common examples of queues. Our computer laboratory has 30 computers networked with a single printer. When students want to print, their print tasks “get in line” with all the other printing tasks that are waiting. The first task in is the next to be completed. If you are last in line, you must wait for all the other tasks to print ahead of you. We will explore this interesting example in more detail later.

In addition to printing queues, operating systems use a number of different queues to control processes within a computer. The scheduling of what gets done next is typically based on a queuing algorithm that tries to execute programs as quickly as possible and serve as many users as it can. Also, as we type, sometimes keystrokes get ahead of the characters that appear on the screen. This is due to the computer doing other work at that moment. The keystrokes are being placed in a queue-like buffer so that they can eventually be displayed on the screen in the proper order.

### 3.4.2. 抽象数据类型 Queue

The queue abstract data type is defined by the following structure and operations. A queue is structured, as described above, as an ordered collection of items which are added at one end, called the “rear,” and removed from the other end, called the “front.” Queues maintain a FIFO ordering property. The queue operations are given below.

- Queue() creates a new queue that is empty. It needs no parameters and returns an empty queue.
- enqueue(item) adds a new item to the rear of the queue. It needs the item and returns nothing.
- dequeue() removes the front item from the queue. It needs no parameters and returns the item. The queue is modified.
- isEmpty() tests to see whether the queue is empty. It needs no parameters and returns a boolean value.
- size() returns the number of items in the queue. It needs no parameters and returns an integer.

As an example, if we assume that q is a queue that has been created and is currently empty, then Table 1 shows the results of a sequence of queue operations. The queue contents are shown such that the front is on the right. 4 was the first item enqueued so it is the first item returned by dequeue.

Queue Operation	Queue Contents	Return Value
q=Queue()	[]	Queue object
q.isEmpty()	[]	True
q.enqueue(4)	[4]	
q.enqueue('dog')	['dog',4]	
q.enqueue(True)	[True,'dog',4]	
q.size()	[True,'dog',4]	3
q.isEmpty()	[True,'dog',4]	False
q.enqueue(8.4)	[8.4,True,'dog',4]	
q.dequeue()	[8.4,True,'dog']	4
q.dequeue()	[8.4,True]	'dog'
q.size()	[8.4,True]	2

表格 17 Example Queue Operations

### 3.4.3. 在 Python 中实现 Queue

It is again appropriate to create a new class for the implementation of the abstract data type queue. As before, we will use the power and simplicity of the list collection to build the internal

representation of the queue.

We need to decide which end of the list to use as the rear and which to use as the front. The implementation shown in Listing 1 assumes that the rear is at position 0 in the list. This allows us to use the insert function on lists to add new elements to the rear of the queue. The pop operation can be used to remove the front element (the last element of the list). Recall that this also means that enqueue will be  $O(n)$  and dequeue will be  $O(1)$ .

Listing 1

```
class Queue:
    def __init__(self):
        self.items = []

    def isEmpty(self):
        return self.items == []

    def enqueue(self, item):
        self.items.insert(0,item)

    def dequeue(self):
        return self.items.pop()

    def size(self):
        return len(self.items)
```

CodeLens 1 shows the Queue class in action as we perform the sequence of operations from Table 1.

```
1  class Queue:
2      def __init__(self):
3          self.items = []
4
5      def isEmpty(self):
6          return self.items == []
7
8      def enqueue(self, item):
9          self.items.insert(0,item)
10
11     def dequeue(self):
12         return self.items.pop()
13
14     def size(self):
15         return len(self.items)
```

```
16
17 q=Queue()
18
19 q.enqueue(4)
20 q.enqueue('dog')
21 q.enqueue(True)
22 print(q.size())
```

Further manipulation of this queue would give the following results:

```
>>> q.size()
3
>>> q.isEmpty()
False
>>> q.enqueue(8.4)
>>> q.dequeue()
4
>>> q.dequeue()
'dog'
>>> q.size()
2
```

### Self Check

Q-9: Suppose you have the following series of queue operations.

```
q = Queue()
q.enqueue('hello')
q.enqueue('dog')
q.enqueue(3)
q.dequeue()
```

What items are left on the queue?

- a) 'hello', 'dog'
- b) 'dog', 3
- c) 'hello', 3
- d) 'hello', 'dog', 3

### 3.4.4. 模拟算法：热土豆

One of the typical applications for showing a queue in action is to simulate a real situation that requires data to be managed in a FIFO manner. To begin, let's consider the children's game Hot Potato. In this game (see Figure 2) children line up in a circle and pass an item from neighbor to neighbor as fast as they can. At a certain point in the game, the action is stopped and the child who



has the item (the potato) is removed from the circle. Play continues until only one child is left.

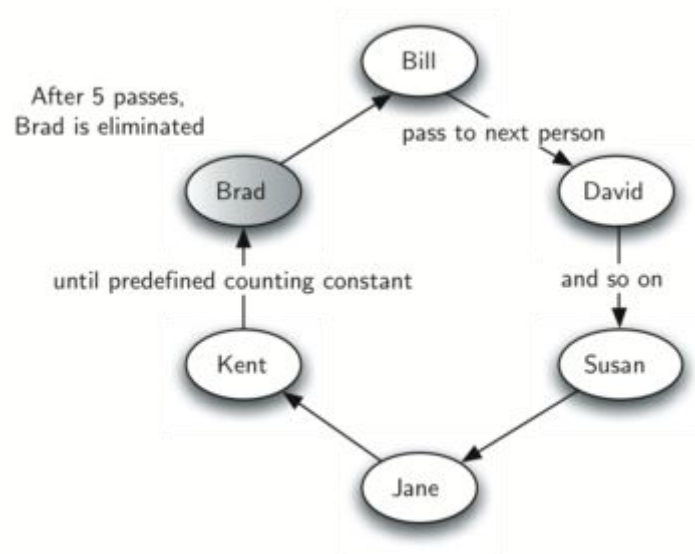


图 29 A Six Person Game of Hot Potato

This game is a modern-day equivalent of the famous Josephus problem. Based on a legend about the famous first-century historian Flavius Josephus, the story is told that in the Jewish revolt against Rome, Josephus and 39 of his comrades held out against the Romans in a cave. With defeat imminent, they decided that they would rather die than be slaves to the Romans. They arranged themselves in a circle. One man was designated as number one, and proceeding clockwise they killed every seventh man. Josephus, according to the legend, was among other things an accomplished mathematician. He instantly figured out where he ought to sit in order to be the last to go. When the time came, instead of killing himself, he joined the Roman side. You can find many different versions of this story. Some count every third man and some allow the last man to escape on a horse. In any case, the idea is the same.

We will implement a general simulation of Hot Potato. Our program will input a list of names and a constant, call it “num,” to be used for counting. It will return the name of the last person remaining after repetitive counting by num. What happens at that point is up to you.

To simulate the circle, we will use a queue (see Figure 3). Assume that the child holding the potato will be at the front of the queue. Upon passing the potato, the simulation will simply dequeue and then immediately enqueue that child, putting her at the end of the line. She will then wait until all the others have been at the front before it will be her turn again. After num dequeue/enqueue operations, the child at the front will be removed permanently and another cycle will begin. This process will continue until only one name remains (the size of the queue is 1).

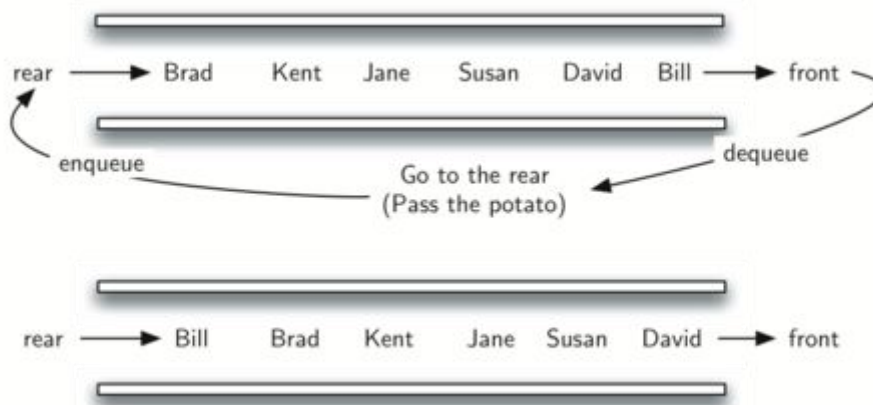


图 30 A Queue Implementation of Hot Potato

The program is shown in ActiveCode 1. A call to the hotPotato function using 7 as the counting constant returns Susan.

```
from pythonds.basic.queue import Queue

def hotPotato(namelist, num):
    simqueue = Queue()
    for name in namelist:
        simqueue.enqueue(name)

    while simqueue.size() > 1:
        for i in range(num):
            simqueue.enqueue(simqueue.dequeue())

        simqueue.dequeue()

    return simqueue.dequeue()

print(hotPotato(["Bill", "David", "Susan", "Jane", "Kent", "Brad"], 7))
```

代码 30 Hot Potato Simulation (qujosephussim)

Note that in this example the value of the counting constant is greater than the number of names in the list. This is not a problem since the queue acts like a circle and counting continues back at the beginning until the value is reached. Also, notice that the list is loaded into the queue such that the first name on the list will be at the front of the queue. Bill in this case is the first item in the list and therefore moves to the front of the queue. A variation of this implementation, described in the exercises, allows for a random counter.

### 3.4.5. 模拟算法：打印任务

A more interesting simulation allows us to study the behavior of the printing queue described earlier in this section. Recall that as students send printing tasks to the shared printer, the tasks are placed in a queue to be processed in a first-come first-served manner. Many questions arise with this configuration. The most important of these might be whether the printer is capable of handling a certain amount of work. If it cannot, students will be waiting too long for printing and may miss their next class.

Consider the following situation in a computer science laboratory. On any average day about 10 students are working in the lab at any given hour. These students typically print up to twice during that time, and the length of these tasks ranges from 1 to 20 pages. The printer in the lab is older, capable of processing 10 pages per minute of draft quality. The printer could be switched to give better quality, but then it would produce only five pages per minute. The slower printing speed could make students wait too long. What page rate should be used?

We could decide by building a simulation that models the laboratory. We will need to construct representations for students, printing tasks, and the printer (Figure 4). As students submit printing tasks, we will add them to a waiting list, a queue of print tasks attached to the printer. When the printer completes a task, it will look at the queue to see if there are any remaining tasks to process. Of interest for us is the average amount of time students will wait for their papers to be printed. This is equal to the average amount of time a task waits in the queue.

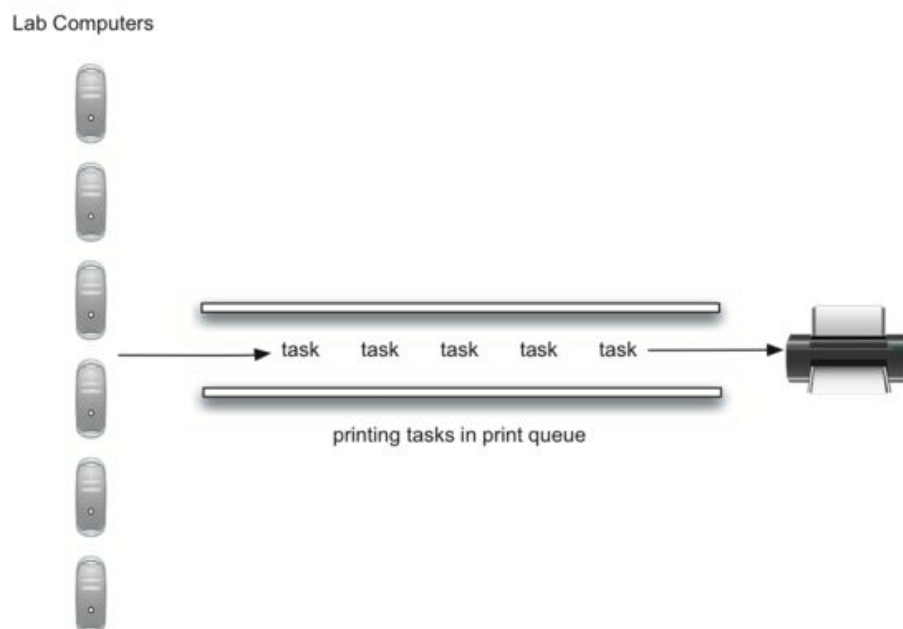


图 31 Computer Science Laboratory Printing Queue

To model this situation we need to use some probabilities. For example, students may print a paper from 1 to 20 pages in length. If each length from 1 to 20 is equally likely, the actual length for a print task can be simulated by using a random number between 1 and 20 inclusive. This means

that there is equal chance of any length from 1 to 20 appearing.

If there are 10 students in the lab and each prints twice, then there are 20 print tasks per hour on average. What is the chance that at any given second, a print task is going to be created? The way to answer this is to consider the ratio of tasks to time. Twenty tasks per hour means that on average there will be one task every 180 seconds:

$$\frac{20 \text{ tasks}}{1 \text{ hour}} \times \frac{1 \text{ hour}}{60 \text{ minutes}} \times \frac{1 \text{ minute}}{60 \text{ seconds}} = \frac{1 \text{ task}}{180 \text{ seconds}}$$

For every second we can simulate the chance that a print task occurs by generating a random number between 1 and 180 inclusive. If the number is 180, we say a task has been created. Note that it is possible that many tasks could be created in a row or we may wait quite a while for a task to appear. That is the nature of simulation. You want to simulate the real situation as closely as possible given that you know general parameters.

### 3.4.6. 主要模拟步骤

Here is the main simulation.

1. Create a queue of print tasks. Each task will be given a timestamp upon its arrival. The queue is empty to start.
2. For each second (currentSecond):
  - a) Does a new print task get created? If so, add it to the queue with the currentSecond as the timestamp.
  - b) If the printer is not busy and if a task is waiting,
    - i. Remove the next task from the print queue and assign it to the printer.
    - ii. Subtract the timestamp from the currentSecond to compute the waiting time for that task.
    - iii. Append the waiting time for that task to a list for later processing.
    - iv. Based on the number of pages in the print task, figure out how much time will be required.
  - c) The printer now does one second of printing if necessary. It also subtracts one second from the time required for that task.
  - d) If the task has been completed, in other words the time required has reached zero, the printer is no longer busy.
3. After the simulation is complete, compute the average waiting time from the list of waiting times generated.

### 3.4.7. Python 实现

To design this simulation we will create classes for the three real-world objects described above: Printer, Task, and PrintQueue.

The Printer class (Listing 2) will need to track whether it has a current task. If it does, then it is busy (lines 13 – 17) and the amount of time needed can be computed from the number of pages in the task. The constructor will also allow the pages-per-minute setting to be initialized. The tick method decrements the internal timer and sets the printer to idle (line 11) if the task is completed.

#### Listing 2

```
class Printer:
    def __init__(self, ppm):
        self.pagerate = ppm
        self.currentTask = None
        self.timeRemaining = 0

    def tick(self):
        if self.currentTask != None:
            self.timeRemaining = self.timeRemaining - 1
            if self.timeRemaining <= 0:
                self.currentTask = None

    def busy(self):
        if self.currentTask != None:
            return True
        else:
            return False

    def startNext(self, newtask):
        self.currentTask = newtask
        self.timeRemaining = newtask.getPages() * 60/self.pagerate
```

The Task class (Listing 3) will represent a single printing task. When the task is created, a random number generator will provide a length from 1 to 20 pages. We have chosen to use the randrange function from the random module.

#### Listing 3

```
>>> import random
>>> random.randrange(1,21)
18
>>> random.randrange(1,21)
8
>>>
```

Each task will also need to keep a timestamp to be used for computing waiting time. This timestamp will represent the time that the task was created and placed in the printer queue. The waitTime method can then be used to retrieve the amount of time spent in the queue before

printing begins.

Listing 3

```
import random

class Task:
    def __init__(self,time):
        self.timestamp = time
        self.pages = random.randrange(1,21)

    def getStamp(self):
        return self.timestamp

    def getPages(self):
        return self.pages

    def waitTime(self, currenttime):
        return currenttime - self.timestamp
```

The main simulation (Listing 4) implements the algorithm described above. The printQueue object is an instance of our existing queue ADT. A boolean helper function, newPrintTask, decides whether a new printing task has been created. We have again chosen to use the randrange function from the random module to return a random integer between 1 and 180. Print tasks arrive once every 180 seconds. By arbitrarily choosing 180 from the range of random integers (line 32), we can simulate this random event. The simulation function allows us to set the total time and the pages per minute for the printer.

Listing 4

```
from pythonds.basic.queue import Queue

import random

def simulation(numSeconds, pagesPerMinute):

    labprinter = Printer(pagesPerMinute)
    printQueue = Queue()
    waitingtimes = []

    for currentSecond in range(numSeconds):

        if newPrintTask():
            task = Task(currentSecond)
```

```

        printQueue.enqueue(task)

    if (not labprinter.busy()) and (not printQueue.isEmpty()):
        nexttask = printQueue.dequeue()
        waitingtimes.append(nexttask.waitTime(currentSecond))
        labprinter.startNext(nexttask)

    labprinter.tick()

    averageWait=sum(waitingtimes)/len(waitingtimes)
    print("Average      Wait      %6.2f      secs      %3d      tasks
remaining."%(averageWait,printQueue.size()))

def newPrintTask():
    num = random.randrange(1,181)
    if num == 180:
        return True
    else:
        return False

for i in range(10):
    simulation(3600,5)

```

When we run the simulation, we should not be concerned that the results are different each time. This is due to the probabilistic nature of the random numbers. We are interested in the trends that may be occurring as the parameters to the simulation are adjusted. Here are some results.

First, we will run the simulation for a period of 60 minutes (3,600 seconds) using a page rate of five pages per minute. In addition, we will run 10 independent trials. Remember that because the simulation works with random numbers each run will return different results.

```

>>>for i in range(10):
    simulation(3600,5)

Average Wait 165.38 secs 2 tasks remaining.
Average Wait 95.07 secs 1 tasks remaining.
Average Wait 65.05 secs 2 tasks remaining.
Average Wait 99.74 secs 1 tasks remaining.
Average Wait 17.27 secs 0 tasks remaining.
Average Wait 239.61 secs 5 tasks remaining.
Average Wait 75.11 secs 1 tasks remaining.
Average Wait 48.33 secs 0 tasks remaining.
Average Wait 39.31 secs 3 tasks remaining.
Average Wait 376.05 secs 1 tasks remaining.

```

After running our 10 trials we can see that the mean average wait time is 122.155 seconds. You can also see that there is a large variation in the average wait time with a minimum average of 17.27 seconds and a maximum of 239.61 seconds. You may also notice that in only two of the cases were all the tasks completed.

Now, we will adjust the page rate to 10 pages per minute, and run the 10 trials again, with a faster page rate our hope would be that more tasks would be completed in the one hour time frame.

```
>>>for i in range(10):
    simulation(3600,10)

Average Wait  1.29 secs 0 tasks remaining.
Average Wait  7.00 secs 0 tasks remaining.
Average Wait 28.96 secs 1 tasks remaining.
Average Wait 13.55 secs 0 tasks remaining.
Average Wait 12.67 secs 0 tasks remaining.
Average Wait  6.46 secs 0 tasks remaining.
Average Wait 22.33 secs 0 tasks remaining.
Average Wait 12.39 secs 0 tasks remaining.
Average Wait  7.27 secs 0 tasks remaining.
Average Wait 18.17 secs 0 tasks remaining.
```

You can run the simulation for yourself in ActiveCode 2.

```
from pythonds.basic.queue import Queue

import random

class Printer:
    def __init__(self, ppm):
        self.pagerate = ppm
        self.currentTask = None
        self.timeRemaining = 0

    def tick(self):
        if self.currentTask != None:
            self.timeRemaining = self.timeRemaining - 1
            if self.timeRemaining <= 0:
                self.currentTask = None

    def busy(self):
        if self.currentTask != None:
            return True
        else:
```



```

        return False

    def startNext(self,newtask):
        self.currentTask = newtask
        self.timeRemaining = newtask.getPages() * 60/self.pagerate

class Task:
    def __init__(self,time):
        self.timestamp = time
        self.pages = random.randrange(1,21)

    def getStamp(self):
        return self.timestamp

    def getPages(self):
        return self.pages

    def waitTime(self, currenttime):
        return currenttime - self.timestamp

def simulation(numSeconds, pagesPerMinute):

    labprinter = Printer(pagesPerMinute)
    printQueue = Queue()
    waitingtimes = []

    for currentSecond in range(numSeconds):

        if newPrintTask():
            task = Task(currentSecond)
            printQueue.enqueue(task)

        if (not labprinter.busy()) and (not printQueue.isEmpty()):
            nexttask = printQueue.dequeue()
            waitingtimes.append( nexttask.waitTime(currentSecond))
            labprinter.startNext(nexttask)

        labprinter.tick()

    averageWait=sum(waitingtimes)/len(waitingtimes)
    print("Average      Wait      %6.2f      secs      %3d      tasks
remaining."%(averageWait,printQueue.size()))

```

```
def newPrintTask():
    num = random.randrange(1,181)
    if num == 180:
        return True
    else:
        return False

for i in range(10):
    simulation(3600,5)
```

代码 31 Printer Queue Simulation (qumainsim)

### 3.4.8. 讨论

We were trying to answer a question about whether the current printer could handle the task load if it were set to print with a better quality but slower page rate. The approach we took was to write a simulation that modeled the printing tasks as random events of various lengths and arrival times.

The output above shows that with 5 pages per minute printing, the average waiting time varied from a low of 17 seconds to a high of 376 seconds (about 6 minutes). With a faster printing rate, the low value was 1 second with a high of only 28. In addition, in 8 out of 10 runs at 5 pages per minute there were print tasks still waiting in the queue at the end of the hour.

Therefore, we are perhaps persuaded that slowing the printer down to get better quality may not be a good idea. Students cannot afford to wait that long for their papers, especially when they need to be getting on to their next class. A six-minute wait would simply be too long.

This type of simulation analysis allows us to answer many questions, commonly known as “what if” questions. All we need to do is vary the parameters used by the simulation and we can simulate any number of interesting behaviors. For example,

- What if enrollment goes up and the average number of students increases by 20?
- What if it is Saturday and students are not needing to get to class? Can they afford to wait?
- What if the size of the average print task decreases since Python is such a powerful language and programs tend to be much shorter?

These questions could all be answered by modifying the above simulation. However, it is important to remember that the simulation is only as good as the assumptions that are used to build it. Real data about the number of print tasks per hour and the number of students per hour was necessary to construct a robust simulation.

#### Self Check

How would you modify the printer simulation to reflect a larger number of students? Suppose that the number of students was doubled. You may need to make some reasonable assumptions about

how this simulation was put together but what would you change? Modify the code. Also suppose that the length of the average print task was cut in half. Change the code to reflect that change. Finally How would you parameterize the number of students, rather than changing the code we would like to make the number of students a parameter of the simulation.

## 3.5. 双端队列 deque

### 3.5.1. 什么是双端队列 deque

A deque, also known as a double-ended queue, is an ordered collection of items similar to the queue. It has two ends, a front and a rear, and the items remain positioned in the collection. What makes a deque different is the unrestrictive nature of adding and removing items. New items can be added at either the front or the rear. Likewise, existing items can be removed from either end. In a sense, this hybrid linear structure provides all the capabilities of stacks and queues in a single data structure. Figure 1 shows a deque of Python data objects.

It is important to note that even though the deque can assume many of the characteristics of stacks and queues, it does not require the LIFO and FIFO orderings that are enforced by those data structures. It is up to you to make consistent use of the addition and removal operations.

../\_images/basicdeque.png

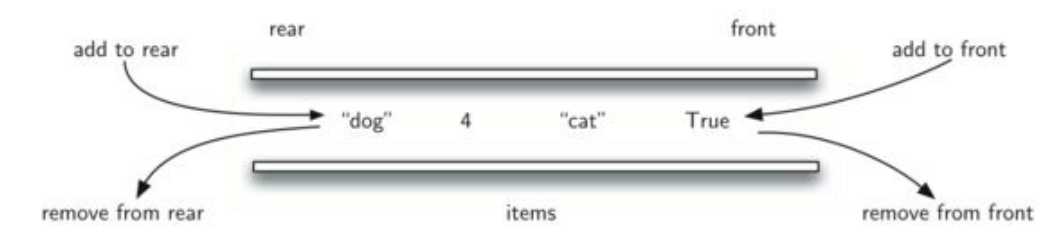


图 32 A Deque of Python Data Objects

### 3.5.2. 抽象数据类型 Deque

The deque abstract data type is defined by the following structure and operations. A deque is structured, as described above, as an ordered collection of items where items are added and removed from either end, either front or rear. The deque operations are given below.

- Deque() creates a new deque that is empty. It needs no parameters and returns an empty deque.
- addFront(item) adds a new item to the front of the deque. It needs the item and returns nothing.
- addRear(item) adds a new item to the rear of the deque. It needs the item and returns nothing.
- removeFront() removes the front item from the deque. It needs no parameters and returns the

item. The deque is modified.

- `removeRear()` removes the rear item from the deque. It needs no parameters and returns the item. The deque is modified.
- `isEmpty()` tests to see whether the deque is empty. It needs no parameters and returns a boolean value.
- `size()` returns the number of items in the deque. It needs no parameters and returns an integer.

As an example, if we assume that `d` is a deque that has been created and is currently empty, then Table {dequeoperations} shows the results of a sequence of deque operations. Note that the contents in front are listed on the right. It is very important to keep track of the front and the rear as you move items in and out of the collection as things can get a bit confusing.

Deque Operation	Deque Contents	Return Value
<code>d=Deque()</code>	<code>[]</code>	Deque object
<code>d.isEmpty()</code>	<code>[]</code>	True
<code>d.addRear(4)</code>	<code>[4]</code>	
<code>d.addRear('dog')</code>	<code>['dog',4,]</code>	
<code>d.addFront('cat')</code>	<code>['dog',4,'cat']</code>	
<code>d.addFront(True)</code>	<code>['dog',4,'cat',True]</code>	
<code>d.size()</code>	<code>['dog',4,'cat',True]</code>	4
<code>d.isEmpty()</code>	<code>['dog',4,'cat',True]</code>	False
<code>d.addRear(8.4)</code>	<code>[8.4,'dog',4,'cat',True]</code>	
<code>d.removeRear()</code>	<code>['dog',4,'cat',True]</code>	8.4
<code>d.removeFront()</code>	<code>['dog',4,'cat']</code>	True

表格 18 Examples of Deque Operations

### 3.5.3. 在 Python 中实现 Deque

As we have done in previous sections, we will create a new class for the implementation of the abstract data type deque. Again, the Python list will provide a very nice set of methods upon which to build the details of the deque. Our implementation (Listing 1) will assume that the rear of the deque is at position 0 in the list.

Listing 1

```
class Deque:
    def __init__(self):
        self.items = []

    def isEmpty(self):
        return self.items == []
```

```
def addFront(self, item):
    self.items.append(item)

def addRear(self, item):
    self.items.insert(0,item)

def removeFront(self):
    return self.items.pop()

def removeRear(self):
    return self.items.pop(0)

def size(self):
    return len(self.items)
```

In `removeFront` we use the `pop` method to remove the last element from the list. However, in `removeRear`, the `pop(0)` method must remove the first element of the list. Likewise, we need to use the `insert` method (line 12) in `addRear` since the `append` method assumes the addition of a new element to the end of the list.

CodeLens 1 shows the `Deque` class in action as we perform the sequence of operations from Table 1.

```
1  class Deque:
2      def __init__(self):
3          self.items = []
4
5      def isEmpty(self):
6          return self.items == []
7
8      def addFront(self, item):
9          self.items.append(item)
10
11     def addRear(self, item):
12         self.items.insert(0,item)
13
14     def removeFront(self):
15         return self.items.pop()
16
17     def removeRear(self):
18         return self.items.pop(0)
19
20     def size(self):
21         return len(self.items)
```

```
22
23 d=Deque()
24 print(d.isEmpty())
25 d.addRear(4)
26 d.addRear('dog')
27 d.addFront('cat')
28 d.addFront(True)
29 print(d.size())
30 print(d.isEmpty())
31 d.addRear(8.4)
32 print(d.removeRear())
33 print(d.removeFront())
```

代码 32 Example Deque Operations (deqtest)

### 3.5.4. “回文词” 判定

An interesting problem that can be easily solved using the deque data structure is the classic palindrome problem. A palindrome is a string that reads the same forward and backward, for example, radar, toot, and madam. We would like to construct an algorithm to input a string of characters and check whether it is a palindrome.

The solution to this problem will use a deque to store the characters of the string. We will process the string from left to right and add each character to the rear of the deque. At this point, the deque will be acting very much like an ordinary queue. However, we can now make use of the dual functionality of the deque. The front of the deque will hold the first character of the string and the rear of the deque will hold the last character (see Figure 2).

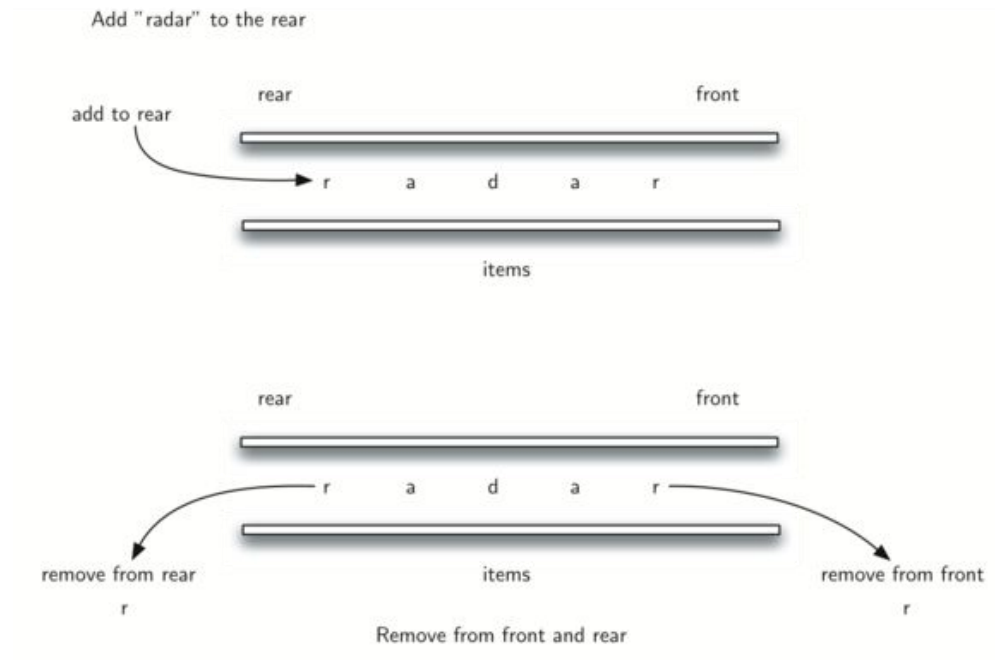


图 33 A Deque

Since we can remove both of them directly, we can compare them and continue only if they match. If we can keep matching first and the last items, we will eventually either run out of characters or be left with a deque of size 1 depending on whether the length of the original string was even or odd. In either case, the string must be a palindrome. The complete function for palindrome-checking appears in ActiveCode 1.

```
from pythonds.basic.deque import Deque

def palchecker(aString):
    chardeque = Deque()

    for ch in aString:
        chardeque.addRear(ch)

    stillEqual = True

    while chardeque.size() > 1 and stillEqual:
        first = chardeque.removeFront()
        last = chardeque.removeRear()
        if first != last:
            stillEqual = False

    return stillEqual

print(palchecker("lsdkjfskf"))
```

```
print(palchecker("radar"))
```

代码 33 A Palindrome Checker Using Deque (palchecker)

## 3.6. 列表 List

Throughout the discussion of basic data structures, we have used Python lists to implement the abstract data types presented. The list is a powerful, yet simple, collection mechanism that provides the programmer with a wide variety of operations. However, not all programming languages include a list collection. In these cases, the notion of a list must be implemented by the programmer.

A list is a collection of items where each item holds a relative position with respect to the others. More specifically, we will refer to this type of list as an unordered list. We can consider the list as having a first item, a second item, a third item, and so on. We can also refer to the beginning of the list (the first item) or the end of the list (the last item). For simplicity we will assume that lists cannot contain duplicate items.

For example, the collection of integers 54, 26, 93, 17, 77, and 31 might represent a simple unordered list of exam scores. Note that we have written them as comma-delimited values, a common way of showing the list structure. Of course, Python would show this list as [54,26,93,17,77,31].

### 3.6.1.抽象数据类型无序列表 Unordered List

The structure of an unordered list, as described above, is a collection of items where each item holds a relative position with respect to the others. Some possible unordered list operations are given below.

- List() creates a new list that is empty. It needs no parameters and returns an empty list.
- add(item) adds a new item to the list. It needs the item and returns nothing. Assume the item is not already in the list.
- remove(item) removes the item from the list. It needs the item and modifies the list. Assume the item is present in the list.
- search(item) searches for the item in the list. It needs the item and returns a boolean value.
- isEmpty() tests to see whether the list is empty. It needs no parameters and returns a boolean value.
- size() returns the number of items in the list. It needs no parameters and returns an integer.
- append(item) adds a new item to the end of the list making it the last item in the collection. It needs the item and returns nothing. Assume the item is not already in the list.
- index(item) returns the position of item in the list. It needs the item and returns the index.



Assume the item is in the list.

- `insert(pos,item)` adds a new item to the list at position `pos`. It needs the item and returns nothing. Assume the item is not already in the list and there are enough existing items to have position `pos`.
- `pop()` removes and returns the last item in the list. It needs nothing and returns an item. Assume the list has at least one item.
- `pop(pos)` removes and returns the item at position `pos`. It needs the position and returns the item. Assume the item is in the list.

### 3.6.2.采用链表实现无序列表

In order to implement an unordered list, we will construct what is commonly known as a linked list. Recall that we need to be sure that we can maintain the relative positioning of the items. However, there is no requirement that we maintain that positioning in contiguous memory. For example, consider the collection of items shown in Figure 1. It appears that these values have been placed randomly. If we can maintain some explicit information in each item, namely the location of the next item (see Figure 2), then the relative position of each item can be expressed by simply following the link from one item to the next.



图 34 Items Not Constrained in Their Physical Placement

../\_images/idea2.png

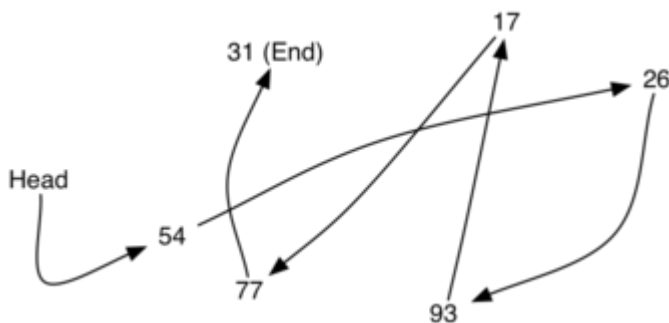


图 35 Relative Positions Maintained by Explicit Links

It is important to note that the location of the first item of the list must be explicitly specified. Once we know where the first item is, the first item can tell us where the second is, and so on. The external reference is often referred to as the head of the list. Similarly, the last item needs to know that there is no next item.

### 3.6.2.1. 类：节点 Node

The basic building block for the linked list implementation is the node. Each node object must hold at least two pieces of information. First, the node must contain the list item itself. We will call this the data field of the node. In addition, each node must hold a reference to the next node. Listing 1 shows the Python implementation. To construct a node, you need to supply the initial data value for the node. Evaluating the assignment statement below will yield a node object containing the value 93 (see Figure 3). You should note that we will typically represent a node object as shown in Figure 4. The Node class also includes the usual methods to access and modify the data and the next reference.

Listing 1

```
class Node:
    def __init__(self,initdata):
        self.data = initdata
        self.next = None

    def getData(self):
        return self.data

    def getNext(self):
        return self.next

    def setData(self,newdata):
        self.data = newdata

    def setNext(self,newnext):
        self.next = newnext
```

We create Node objects in the usual way.

```
>>> temp = Node(93)
>>> temp.getData()
93
```

The special Python reference value None will play an important role in the Node class and later in the linked list itself. A reference to None will denote the fact that there is no next node. Note in the constructor that a node is initially created with next set to None. Since this is sometimes referred to as “grounding the node,” we will use the standard ground symbol to denote a reference that is referring to None. It is always a good idea to explicitly assign None to your initial next reference values.

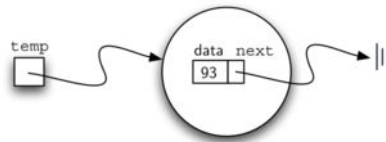


图 36 A Node Object Contains the Item and a Reference to the Next Node



图 37 A Typical Representation for a Node

### 3.6.2.2. 类：无序列表 Unordered List

As we suggested above, the unordered list will be built from a collection of nodes, each linked to the next by explicit references. As long as we know where to find the first node (containing the first item), each item after that can be found by successively following the next links. With this in mind, the UnorderedList class must maintain a reference to the first node. Listing 2 shows the constructor. Note that each list object will maintain a single reference to the head of the list.

#### Listing 2

```
class UnorderedList:

    def __init__(self):
        self.head = None
```

Initially when we construct a list, there are no items. The assignment statement

```
>>> mylist = UnorderedList()
```

creates the linked list representation shown in Figure 5. As we discussed in the Node class, the special reference None will again be used to state that the head of the list does not refer to anything. Eventually, the example list given earlier will be represented by a linked list as shown in Figure 6. The head of the list refers to the first node which contains the first item of the list. In turn, that node holds a reference to the next node (the next item) and so on. It is very important to note that the list class itself does not contain any node objects. Instead it contains a single reference to only the first node in the linked structure.

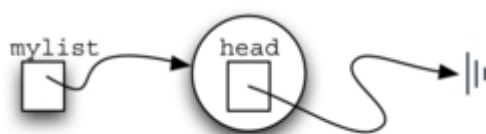


Figure 5: An Empty List

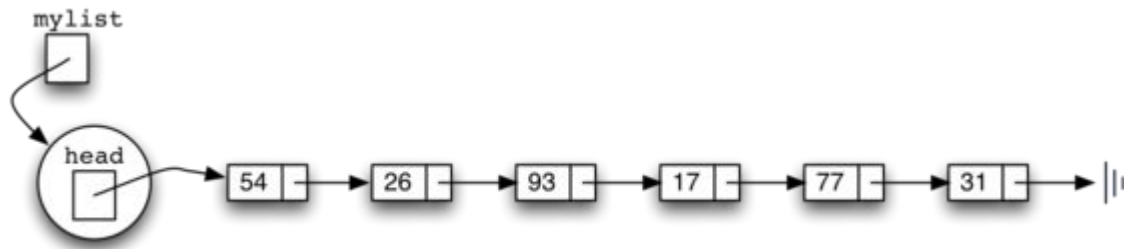


Figure 6: A Linked List of Integers

The `isEmpty` method, shown in Listing 3, simply checks to see if the head of the list is a reference to `None`. The result of the boolean expression `self.head==None` will only be true if there are no nodes in the linked list. Since a new list is empty, the constructor and the check for empty must be consistent with one another. This shows the advantage to using the reference `None` to denote the “end” of the linked structure. In Python, `None` can be compared to any reference. Two references are equal if they both refer to the same object. We will use this often in our remaining methods.

Listing 3

```
def isEmpty(self):
    return self.head == None
```

So, how do we get items into our list? We need to implement the `add` method. However, before we can do that, we need to address the important question of where in the linked list to place the new item. Since this list is unordered, the specific location of the new item with respect to the other items already in the list is not important. The new item can go anywhere. With that in mind, it makes sense to place the new item in the easiest location possible.

Recall that the linked list structure provides us with only one entry point, the head of the list. All of the other nodes can only be reached by accessing the first node and then following next links. This means that the easiest place to add the new node is right at the head, or beginning, of the list. In other words, we will make the new item the first item of the list and the existing items will need to be linked to this new first item so that they follow.

The linked list shown in Figure 6 was built by calling the `add` method a number of times.

```
>>> mylist.add(31)
>>> mylist.add(77)
>>> mylist.add(17)
>>> mylist.add(93)
>>> mylist.add(26)
>>> mylist.add(54)
```

Note that since 31 is the first item added to the list, it will eventually be the last node on the linked list as every other item is added ahead of it. Also, since 54 is the last item added, it will become

the data value in the first node of the linked list.

The add method is shown in Listing 4. Each item of the list must reside in a node object. Line 2 creates a new node and places the item as its data. Now we must complete the process by linking the new node into the existing structure. This requires two steps as shown in Figure 7. Step 1 (line 3) changes the next reference of the new node to refer to the old first node of the list. Now that the rest of the list has been properly attached to the new node, we can modify the head of the list to refer to the new node. The assignment statement in line 4 sets the head of the list.

The order of the two steps described above is very important. What happens if the order of line 3 and line 4 is reversed? If the modification of the head of the list happens first, the result can be seen in Figure 8. Since the head was the only external reference to the list nodes, all of the original nodes are lost and can no longer be accessed.

Listing 4

```
def add(self,item):  
    temp = Node(item)  
    temp.setNext(self.head)  
    self.head = temp
```

../\_images/addtohead.png

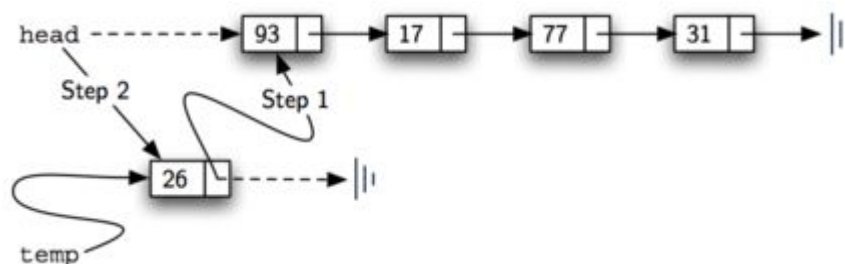


Figure 7: Adding a New Node is a Two-Step Process

../\_images/wrongorder.png

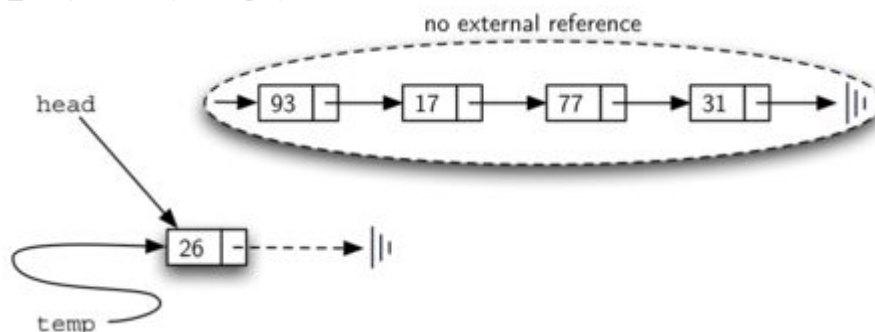


Figure 8: Result of Reversing the Order of the Two Steps

The next methods that we will implement – size, search, and remove – are all based on a technique known as linked list traversal. Traversal refers to the process of systematically visiting each node. To do this we use an external reference that starts at the first node in the list. As we visit each node,

we move the reference to the next node by “traversing” the next reference.

To implement the size method, we need to traverse the linked list and keep a count of the number of nodes that occurred. Listing 5 shows the Python code for counting the number of nodes in the list. The external reference is called `current` and is initialized to the head of the list in line 2. At the start of the process we have not seen any nodes so the count is set to 0. Lines 4 – 6 actually implement the traversal. As long as the current reference has not seen the end of the list (`None`), we move `current` along to the next node via the assignment statement in line 6. Again, the ability to compare a reference to `None` is very useful. Every time `current` moves to a new node, we add 1 to count. Finally, count gets returned after the iteration stops. Figure 9 shows this process as it proceeds down the list.

Listing 5

```
def size(self):
    current = self.head
    count = 0
    while current != None:
        count = count + 1
        current = current.getNext()

    return count
```

../\_images/traversal.png

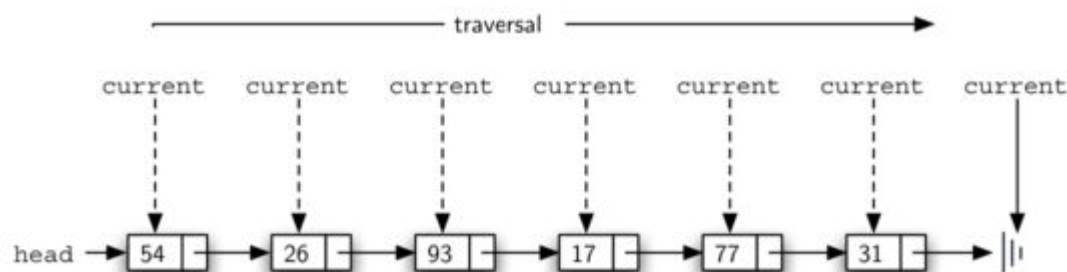


Figure 9: Traversing the Linked List from the Head to the End

Searching for a value in a linked list implementation of an unordered list also uses the traversal technique. As we visit each node in the linked list we will ask whether the data stored there matches the item we are looking for. In this case, however, we may not have to traverse all the way to the end of the list. In fact, if we do get to the end of the list, that means that the item we are looking for must not be present. Also, if we do find the item, there is no need to continue.

Listing 6 shows the implementation for the search method. As in the size method, the traversal is initialized to start at the head of the list (line 2). We also use a boolean variable called `found` to remember whether we have located the item we are searching for. Since we have not found the item at the start of the traversal, `found` can be set to `False` (line 3). The iteration in line 4 takes into account both conditions discussed above. As long as there are more nodes to visit and we have not

found the item we are looking for, we continue to check the next node. The question in line 5 asks whether the data item is present in the current node. If so, found can be set to True.

#### Listing 6

```
def search(self,item):
    current = self.head
    found = False
    while current != None and not found:
        if current.getData() == item:
            found = True
        else:
            current = current.getNext()

    return found
```

As an example, consider invoking the search method looking for the item 17.

```
>>> mylist.search(17)
True
```

Since 17 is in the list, the traversal process needs to move only to the node containing 17. At that point, the variable found is set to True and the while condition will fail, leading to the return value seen above. This process can be seen in Figure 10.

../\_images/search.png

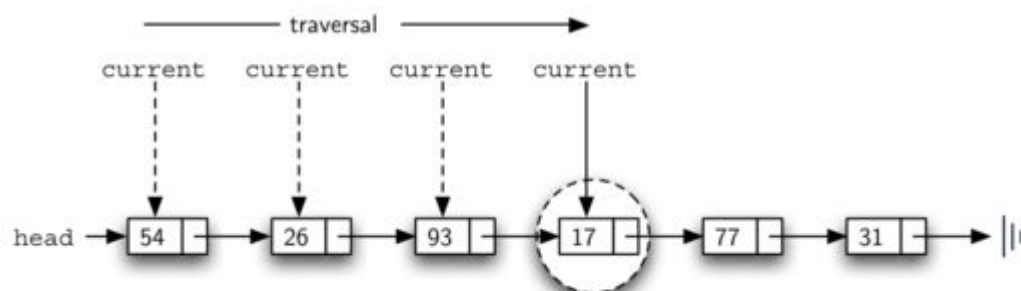


Figure 10: Successful Search for the Value 17

The remove method requires two logical steps. First, we need to traverse the list looking for the item we want to remove. Once we find the item (recall that we assume it is present), we must remove it. The first step is very similar to search. Starting with an external reference set to the head of the list, we traverse the links until we discover the item we are looking for. Since we assume that item is present, we know that the iteration will stop before current gets to None. This means that we can simply use the boolean found in the condition.

When found becomes True, current will be a reference to the node containing the item to be removed. But how do we remove it? One possibility would be to replace the value of the item with

some marker that suggests that the item is no longer present. The problem with this approach is the number of nodes will no longer match the number of items. It would be much better to remove the item by removing the entire node.

In order to remove the node containing the item, we need to modify the link in the previous node so that it refers to the node that comes after current. Unfortunately, there is no way to go backward in the linked list. Since current refers to the node ahead of the node where we would like to make the change, it is too late to make the necessary modification.

The solution to this dilemma is to use two external references as we traverse down the linked list. current will behave just as it did before, marking the current location of the traverse. The new reference, which we will call previous, will always travel one node behind current. That way, when current stops at the node to be removed, previous will be referring to the proper place in the linked list for the modification.

Listing 7 shows the complete remove method. Lines 2 – 3 assign initial values to the two references. Note that current starts out at the list head as in the other traversal examples. previous, however, is assumed to always travel one node behind current. For this reason, previous starts out with a value of None since there is no node before the head (see Figure 11). The boolean variable found will again be used to control the iteration.

In lines 6 – 7 we ask whether the item stored in the current node is the item we wish to remove. If so, found can be set to True. If we do not find the item, previous and current must both be moved one node ahead. Again, the order of these two statements is crucial. previous must first be moved one node ahead to the location of current. At that point, current can be moved. This process is often referred to as “inch-worming” as previous must catch up to current before current moves ahead. Figure 12 shows the movement of previous and current as they progress down the list looking for the node containing the value 17.

Listing 7

```
def remove(self,item):
    current = self.head
    previous = None
    found = False
    while not found:
        if current.getData() == item:
            found = True
        else:
            previous = current
            current = current.getNext()

    if previous == None:
        self.head = current.getNext()
```



```
else:
    previous.setNext(current.getNext())
```

../\_images/removeinit.png

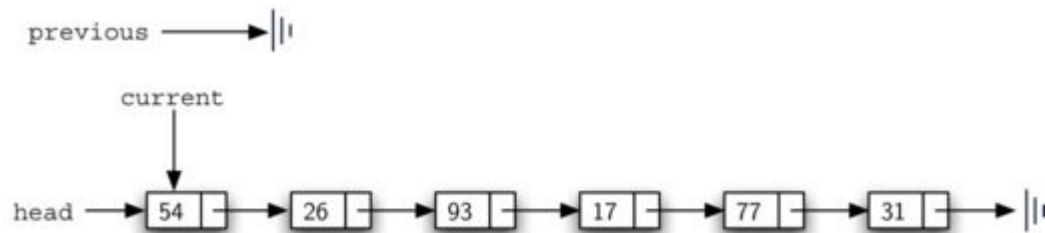


Figure 11: Initial Values for the previous and current References

../\_images/prevcurr.png

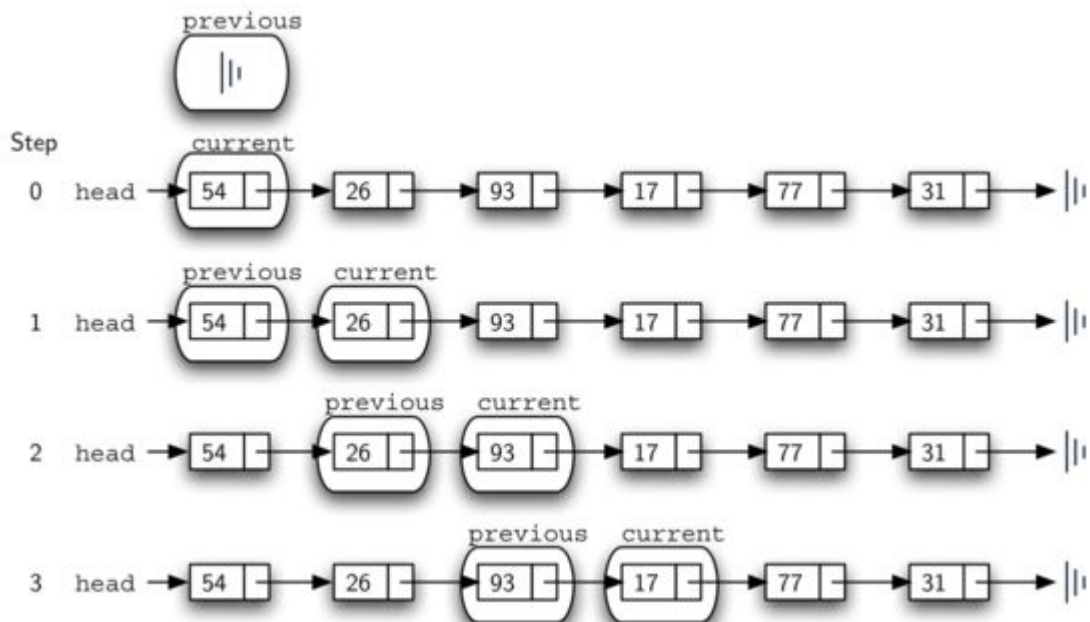


Figure 12: previous and current Move Down the List

Once the searching step of the remove has been completed, we need to remove the node from the linked list. Figure 13 shows the link that must be modified. However, there is a special case that needs to be addressed. If the item to be removed happens to be the first item in the list, then current will reference the first node in the linked list. This also means that previous will be None. We said earlier that previous would be referring to the node whose next reference needs to be modified in order to complete the remove. In this case, it is not previous but rather the head of the list that needs to be changed (see Figure 14).

../\_images/remove.png

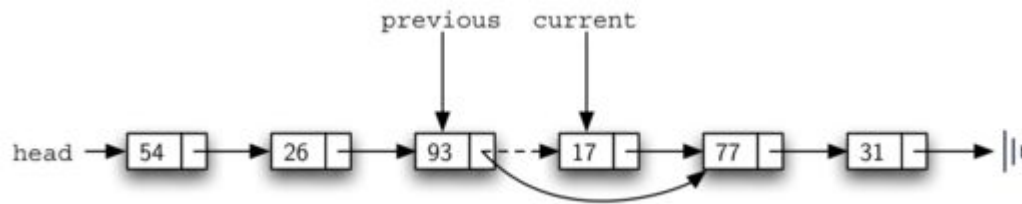


Figure 13: Removing an Item from the Middle of the List

../\_images/remove2.png

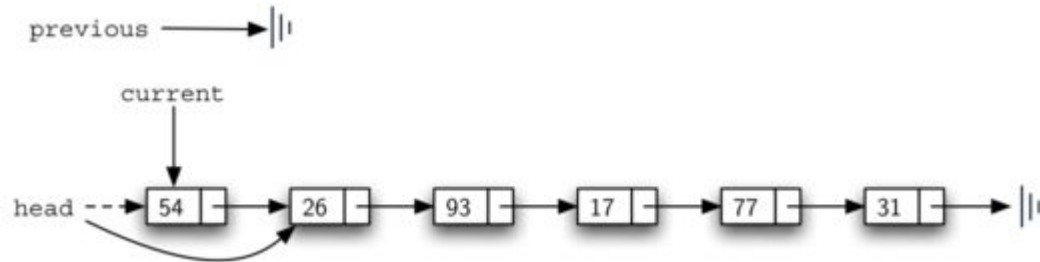


Figure 14: Removing the First Node from the List

Line 12 allows us to check whether we are dealing with the special case described above. If `previous` did not move, it will still have the value `None` when the boolean `found` becomes `True`. In that case (line 13) the `head` of the list is modified to refer to the node after the `current` node, in effect removing the first node from the linked list. However, if `previous` is not `None`, the node to be removed is somewhere down the linked list structure. In this case the `previous` reference is providing us with the node whose `next` reference must be changed. Line 15 uses the `setNext` method from `previous` to accomplish the removal. Note that in both cases the destination of the reference change is `current.getNext()`. One question that often arises is whether the two cases shown here will also handle the situation where the item to be removed is in the last node of the linked list. We leave that for you to consider.

You can try out the `UnorderedList` class in `ActiveCode 1`.

Run Show/Hide Code

The Complete `UnorderedList` Class (`unorderedlistcomplete`)

The remaining methods `append`, `insert`, `index`, and `pop` are left as exercises. Remember that each of these must take into account whether the change is taking place at the head of the list or someplace else. Also, `insert`, `index`, and `pop` require that we name the positions of the list. We will assume that position names are integers starting with 0.

### Self Check

Part I: Implement the `append` method for `UnorderedList`. What is the time complexity of the method you created?

Part II: In the previous problem, you most likely created an `append` method that was  $O(n)$ . If you

add an instance variable to the `UnorderedList` class you can create an `append` method that is  $O(1)$ . Modify your `append` method to be  $O(1)$  Be Careful! To really do this correctly you will need to consider a couple of special cases that may require you to make a modification to the `add` method as well.

### 3.6.3.抽象数据类型：有序列表 **Ordered List**

We will now consider a type of list known as an ordered list. For example, if the list of integers shown above were an ordered list (ascending order), then it could be written as 17, 26, 31, 54, 77, and 93. Since 17 is the smallest item, it occupies the first position in the list. Likewise, since 93 is the largest, it occupies the last position.

The structure of an ordered list is a collection of items where each item holds a relative position that is based upon some underlying characteristic of the item. The ordering is typically either ascending or descending and we assume that list items have a meaningful comparison operation that is already defined. Many of the ordered list operations are the same as those of the unordered list.

- `OrderedList()` creates a new ordered list that is empty. It needs no parameters and returns an empty list.
- `add(item)` adds a new item to the list making sure that the order is preserved. It needs the item and returns nothing. Assume the item is not already in the list.
- `remove(item)` removes the item from the list. It needs the item and modifies the list. Assume the item is present in the list.
- `search(item)` searches for the item in the list. It needs the item and returns a boolean value.
- `isEmpty()` tests to see whether the list is empty. It needs no parameters and returns a boolean value.
- `size()` returns the number of items in the list. It needs no parameters and returns an integer.
- `index(item)` returns the position of item in the list. It needs the item and returns the index. Assume the item is in the list.
- `pop()` removes and returns the last item in the list. It needs nothing and returns an item. Assume the list has at least one item.
- `pop(pos)` removes and returns the item at position `pos`. It needs the position and returns the item. Assume the item is in the list.

### 3.6.4. 实现有序列表

In order to implement the ordered list, we must remember that the relative positions of the items are based on some underlying characteristic. The ordered list of integers given above (17, 26, 31, 54, 77, and 93) can be represented by a linked structure as shown in Figure 15. Again, the node and link structure is ideal for representing the relative positioning of the items.

../\_images/orderlinkedlist.png



Figure 15: An Ordered Linked List

To implement the `OrderedList` class, we will use the same technique as seen previously with unordered lists. Once again, an empty list will be denoted by a head reference to `None` (see Listing 8).

Listing 8

```
class OrderedList:
    def __init__(self):
        self.head = None
```

As we consider the operations for the ordered list, we should note that the `isEmpty` and `size` methods can be implemented the same as with unordered lists since they deal only with the number of nodes in the list without regard to the actual item values. Likewise, the `remove` method will work just fine since we still need to find the item and then link around the node to remove it. The two remaining methods, `search` and `add`, will require some modification.

The search of an unordered linked list required that we traverse the nodes one at a time until we either find the item we are looking for or run out of nodes (`None`). It turns out that the same approach would actually work with the ordered list and in fact in the case where we find the item it is exactly what we need. However, in the case where the item is not in the list, we can take advantage of the ordering to stop the search as soon as possible.

For example, Figure 16 shows the ordered linked list as a search is looking for the value 45. As we traverse, starting at the head of the list, we first compare against 17. Since 17 is not the item we are looking for, we move to the next node, in this case 26. Again, this is not what we want, so we move on to 31 and then on to 54. Now, at this point, something is different. Since 54 is not the item we are looking for, our former strategy would be to move forward. However, due to the fact that this is an ordered list, that will not be necessary. Once the value in the node becomes greater than the item we are searching for, the search can stop and return `False`. There is no way the item could exist further out in the linked list.

../\_images/orderedsearch.png

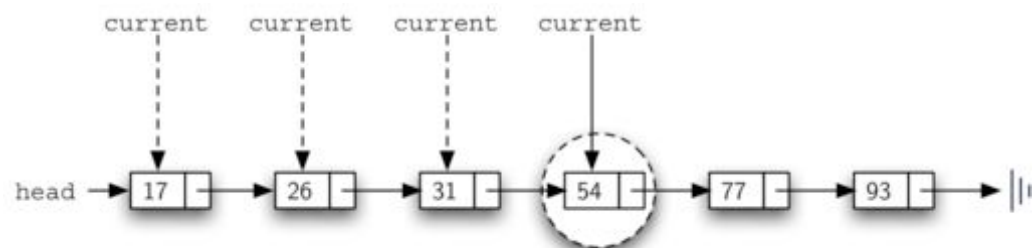


Figure 16: Searching an Ordered Linked List

Listing 9 shows the complete search method. It is easy to incorporate the new condition discussed above by adding another boolean variable, stop, and initializing it to False (line 4). While stop is False (not stop) we can continue to look forward in the list (line 5). If any node is ever discovered that contains data greater than the item we are looking for, we will set stop to True (lines 9 – 10). The remaining lines are identical to the unordered list search.

Listing 9

```
def search(self,item):
    current = self.head
    found = False
    stop = False
    while current != None and not found and not stop:
        if current.getData() == item:
            found = True
        else:
            if current.getData() > item:
                stop = True
            else:
                current = current.getNext()

    return found
```

The most significant method modification will take place in add. Recall that for unordered lists, the add method could simply place a new node at the head of the list. It was the easiest point of access. Unfortunately, this will no longer work with ordered lists. It is now necessary that we discover the specific place where a new item belongs in the existing ordered list.

Assume we have the ordered list consisting of 17, 26, 54, 77, and 93 and we want to add the value 31. The add method must decide that the new item belongs between 26 and 54. Figure 17 shows the setup that we need. As we explained earlier, we need to traverse the linked list looking for the place where the new node will be added. We know we have found that place when either we run out of nodes (current becomes None) or the value of the current node becomes greater than the item we wish to add. In our example, seeing the value 54 causes us to stop.

../\_images/linkedListinsert.png

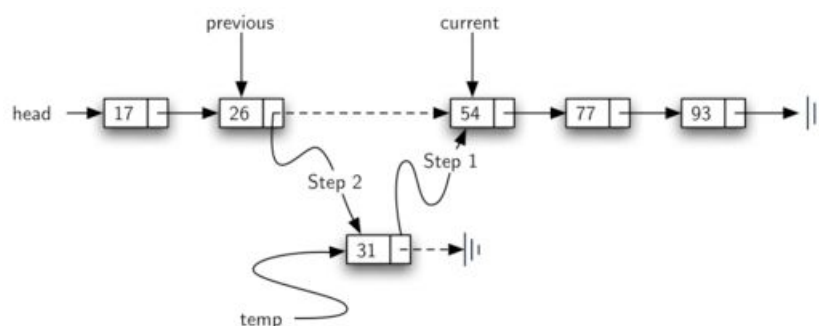


Figure 17: Adding an Item to an Ordered Linked List

As we saw with unordered lists, it is necessary to have an additional reference, again called `previous`, since `current` will not provide access to the node that must be modified. Listing 10 shows the complete `add` method. Lines 2 – 3 set up the two external references and lines 9 – 10 again allow `previous` to follow one node behind `current` every time through the iteration. The condition (line 5) allows the iteration to continue as long as there are more nodes and the value in the current node is not larger than the item. In either case, when the iteration fails, we have found the location for the new node.

The remainder of the method completes the two-step process shown in Figure 17. Once a new node has been created for the item, the only remaining question is whether the new node will be added at the beginning of the linked list or some place in the middle. Again, `previous == None` (line 13) can be used to provide the answer.

Listing 10

```
def add(self,item):
    current = self.head
    previous = None
    stop = False
    while current != None and not stop:
        if current.getData() > item:
            stop = True
        else:
            previous = current
            current = current.getNext()

    temp = Node(item)
    if previous == None:
        temp.setNext(self.head)
        self.head = temp
    else:
        temp.setNext(current)
        previous.setNext(temp)
```

The `OrderedList` class with methods discussed thus far can be found in ActiveCode 1. We leave the remaining methods as exercises. You should carefully consider whether the unordered implementations will work given that the list is now ordered.

Run Show/Hide Code

OrderedList Class Thus Far (orderedlistclass)

### 3.6.5. 链表实现算法分析

#### Analysis of Linked Lists

To analyze the complexity of the linked list operations, we need to consider whether they require traversal. Consider a linked list that has  $n$  nodes. The `isEmpty` method is  $O(1)$  since it requires one step to check the head reference for `None`. `size`, on the other hand, will always require  $n$  steps since there is no way to know how many nodes are in the linked list without traversing from head to end. Therefore, `length` is  $O(n)$ . Adding an item to an unordered list will always be  $O(1)$  since we simply place the new node at the head of the linked list. However, `search` and `remove`, as well as `add` for an ordered list, all require the traversal process. Although on average they may need to traverse only half of the nodes, these methods are all  $O(n)$  since in the worst case each will process every node in the list.

You may also have noticed that the performance of this implementation differs from the actual performance given earlier for Python lists. This suggests that linked lists are not the way Python lists are implemented. The actual implementation of a Python list is based on the notion of an array. We discuss this in more detail in another chapter.

### 3.7. 小结

- Linear data structures maintain their data in an ordered fashion.
- Stacks are simple data structures that maintain a LIFO, last-in first-out, ordering.
- The fundamental operations for a stack are `push`, `pop`, and `isEmpty`.
- Queues are simple data structures that maintain a FIFO, first-in first-out, ordering.
- The fundamental operations for a queue are `enqueue`, `dequeue`, and `isEmpty`.
- Prefix, infix, and postfix are all ways to write expressions.
- Stacks are very useful for designing algorithms to evaluate and translate expressions.
- Stacks can provide a reversal characteristic.
- Queues can assist in the construction of timing simulations.
- Simulations use random number generators to create a real-life situation and allow us to answer “what if” types of questions.
- Deques are data structures that allow hybrid behavior like that of stacks and queues.
- The fundamental operations for a deque are `addFront`, `addRear`, `removeFront`, `removeRear`, and `isEmpty`.
- Lists are collections of items where each item holds a relative position.
- A linked list implementation maintains logical order without requiring physical storage requirements.
- Modification to the head of the linked list is a special case.

### 3.8. 关键词

balanced parentheses	data field	deque
first-in first-out (FIFO)	fully parenthesized	head
infix	last-in first-out (LIFO)	linear data structure
linked list	linked list traversal	list
node	palindrome	postfix
precedence	prefix	queue
simulation	stack	

### 3.9. 问题讨论

- Convert the following values to binary using “divide by 2.” Show the stack of remainders.
  - 17
  - 45
  - 96
- Convert the following infix expressions to prefix (use full parentheses):
  - $(A+B)*(C+D)*(E+F)$
  - $A+((B+C)*(D+E))$
  - $A*B*C*D+E+F$
- Convert the above infix expressions to postfix (use full parentheses).
- Convert the above infix expressions to postfix using the direct conversion algorithm. Show the stack as the conversion takes place.
- Evaluate the following postfix expressions. Show the stack as each operand and operator is processed.
  - $2\ 3\ *\ 4\ +$
  - $1\ 2\ +\ 3\ +\ 4\ +\ 5\ +$
  - $1\ 2\ 3\ 4\ 5\ *\ +\ *\ +$
- The alternative implementation of the Queue ADT is to use a list such that the rear of the queue is at the end of the list. What would this mean for Big-O performance?
- What is the result of carrying out both steps of the linked list add method in reverse order? What kind of reference results? What types of problems may result?
- Explain how the linked list remove method works when the item to be removed is in the last node.
- Explain how the remove method works when the item is in the only node in the linked list.

### 3.10. 编程练习

- Modify the infix-to-postfix algorithm so that it can handle errors.
- Modify the postfix evaluation algorithm so that it can handle errors.
- Implement a direct infix evaluator that combines the functionality of infix-to-postfix conversion and the postfix evaluation algorithm. Your evaluator should process infix tokens



from left to right and use two stacks, one for operators and one for operands, to perform the evaluation.

4. Turn your direct infix evaluator from the previous problem into a calculator.
5. Implement the Queue ADT, using a list such that the rear of the queue is at the end of the list.
6. Design and implement an experiment to do benchmark comparisons of the two queue implementations. What can you learn from such an experiment?
7. It is possible to implement a queue such that both enqueue and dequeue have  $O(1)$  performance on average. In this case it means that most of the time enqueue and dequeue will be  $O(1)$  except in one particular circumstance where dequeue will be  $O(n)$ .
8. Consider a real life situation. Formulate a question and then design a simulation that can help to answer it. Possible situations include:
  - a) Cars lined up at a car wash
  - b) Customers at a grocery store check-out
  - c) Airplanes taking off and landing on a runway
  - d) A bank teller

Be sure to state any assumptions that you make and provide any probabilistic data that must be considered as part of the scenario.

9. Modify the Hot Potato simulation to allow for a randomly chosen counting value so that each pass is not predictable from the previous one.
10. Implement a radix sorting machine. A radix sort for base 10 integers is a mechanical sorting technique that utilizes a collection of bins, one main bin and 10 digit bins. Each bin acts like a queue and maintains its values in the order that they arrive. The algorithm begins by placing each number in the main bin. Then it considers each value digit by digit. The first value is removed and placed in a digit bin corresponding to the digit being considered. For example, if the ones digit is being considered, 534 is placed in digit bin 4 and 667 is placed in digit bin 7. Once all the values are placed in the corresponding digit bins, the values are collected from bin 0 to bin 9 and placed back in the main bin. The process continues with the tens digit, the hundreds, and so on. After the last digit is processed, the main bin contains the values in order.
11. Another example of the parentheses matching problem comes from hypertext markup language (HTML). In HTML, tags exist in both opening and closing forms and must be balanced to properly describe a web document. This very simple HTML document:

```
<html>
  <head>
    <title>
      Example
    </title>
  </head>

  <body>
    <h1>Hello, world</h1>
  </body>
</html>
```

is intended only to show the matching and nesting structure for tags in the language. Write a

program that can check an HTML document for proper opening and closing tags.

12. Extend the program from Listing 2.15 to handle palindromes with spaces. For example, I PREFER PI is a palindrome that reads the same forward and backward if you ignore the blank characters.
13. To implement the length method, we counted the number of nodes in the list. An alternative strategy would be to store the number of nodes in the list as an additional piece of data in the head of the list. Modify the UnorderedList class to include this information and rewrite the length method.
14. Implement the remove method so that it works correctly in the case where the item is not in the list.
15. Modify the list classes to allow duplicates. Which methods will be impacted by this change?
16. Implement the `__str__` method in the UnorderedList class. What would be a good string representation for a list?
17. Implement `__str__` method so that lists are displayed the Python way (with square brackets).
18. Implement the remaining operations defined in the UnorderedList ADT (append, index, pop, insert).
19. Implement a slice method for the UnorderedList class. It should take two parameters, start and stop, and return a copy of the list starting at the start position and going up to but not including the stop position.
20. Implement the remaining operations defined in the OrderedList ADT.
21. Consider the relationship between Unordered and Ordered lists. Is it possible that inheritance could be used to build a more efficient implementation? Implement this inheritance hierarchy.
22. Implement a stack using linked lists.
23. Implement a queue using linked lists.
24. Implement a deque using linked lists.
25. Design and implement an experiment that will compare the performance of a Python list with a list implemented as a linked list.
26. Design and implement an experiment that will compare the performance of the Python list

based stack and queue with the linked list implementation.

27. The linked list implementation given above is called a singly linked list because each node has a single reference to the next node in sequence. An alternative implementation is known as a doubly linked list. In this implementation, each node has a reference to the next node (commonly called next) as well as a reference to the preceding node (commonly called back). The head reference also contains two references, one to the first node in the linked list and one to the last. Code this implementation in Python.
28. Create an implementation of a queue that would have an average

## 4. 递归 Recursion

### 4.1. 目标

- 了解某些难解的问题具有简单的递归解决法；
- 学会如何用递归方式写程序；
- 了解和应用递归的三法则；
- 了解递归是迭代 iteration 的一种形式；
- 实现问题的递归描述；
- 了解递归在计算机系统中是如何实现的。

### 4.2. 什么是递归

**Recursion** is a method of solving problems that involves breaking a problem down into smaller and smaller subproblems until you get to a small enough problem that it can be solved trivially. Usually recursion involves a function calling itself. While it may not seem like much on the surface, recursion allows us to write elegant solutions to problems that may otherwise be very difficult to program.

#### 4.2.1. 计算数列表的和

We will begin our investigation with a simple problem that you already know how to solve without using recursion. Suppose that you want to calculate the sum of a list of numbers such as: [1,3,5,7,9]. An iterative function that computes the sum is shown in ActiveCode 1. The function uses an accumulator variable (theSum) to compute a running total of all the numbers in the list by starting with 0 and adding each number in the list.

```
def listsum(numList):
    theSum = 0
    for i in numList:
        theSum = theSum + i
    return theSum
print(listsum([1,3,5,7,9]))
```

Iterative Summation (lst\_itsum)

Pretend for a minute that you do not have while loops or for loops. How would you compute the sum of a list of numbers? If you were a mathematician you might start by recalling that addition is a function that is defined for two parameters, a pair of numbers. To redefine the problem from adding a list to adding pairs of numbers, we could rewrite the list as a fully parenthesized

expression. Such an expression looks like this:

```
((((1+3)+5)+7)+9)
```

We can also parenthesize the expression the other way around,

```
(1+(3+(5+(7+9))))
```

Notice that the innermost set of parentheses, (7+9), is a problem that we can solve without a loop or any special constructs. In fact, we can use the following sequence of simplifications to compute a final sum.

```
total= (1+(3+(5+(7+9))))  
total= (1+(3+(5+16)))  
total= (1+(3+21))  
total= (1+24)  
total= 25
```

How can we take this idea and turn it into a Python program? First, let's restate the sum problem in terms of Python lists. We might say the the sum of the list numList is the sum of the first element of the list (numList[0]), and the sum of the numbers in the rest of the list (numList[1:]). To state it in a functional form:

```
listSum(numList)=first(numList)+listSum(rest(numList))
```

In this equation first(numList) returns the first element of the list and rest(numList) returns a list of everything but the first element. This is easily expressed in Python as shown in ActiveCode 2.

```
def listsum(numList):  
    if len(numList) == 1:  
        return numList[0]  
    else:  
        return numList[0] + listsum(numList[1:])  
print(listsum([1,3,5,7,9]))
```

#### Recursive Summation (lst\_recsum)

There are a few key ideas in this listing to look at. First, on line 2 we are checking to see if the list is one element long. This check is crucial and is our escape clause from the function. The sum of a list of length 1 is trivial; it is just the number in the list. Second, on line 5 our function calls itself! This is the reason that we call the listsum algorithm recursive. A recursive function is a function that calls itself.

Figure 1 shows the series of recursive calls that are needed to sum the list [1,3,5,7,9]. You should think of this series of calls as a series of simplifications. Each time we make a recursive call we

are solving a smaller problem, until we reach the point where the problem cannot get any smaller.

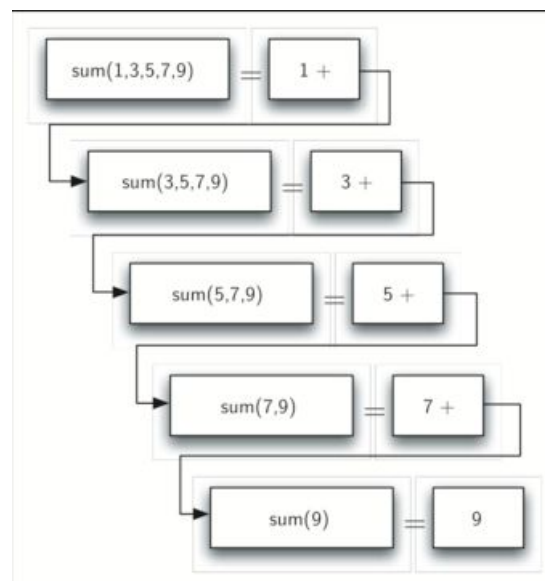


Figure 1: Series of Recursive Calls Adding a List of Numbers

When we reach the point where the problem is as simple as it can get, we begin to piece together the solutions of each of the small problems until the initial problem is solved. Figure 2 shows the additions that are performed as `listsum` works its way backward through the series of calls. When `listsum` returns from the topmost problem, we have the solution to the whole problem.

Image

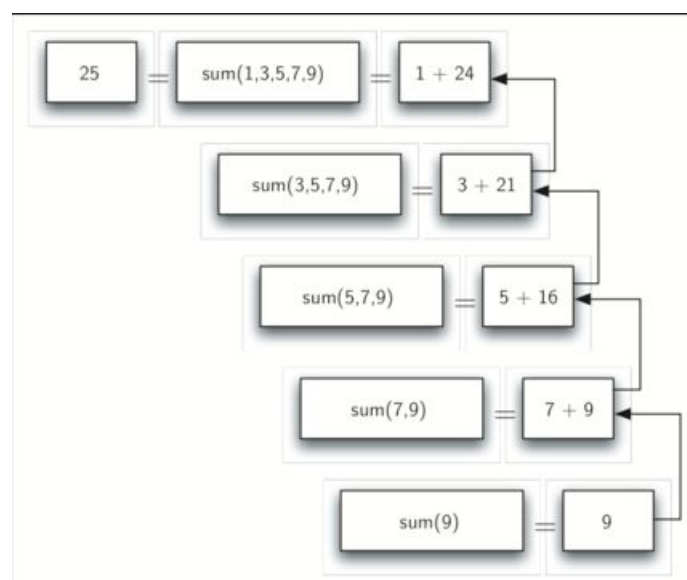


Figure2: Series of Recursive Returns from Adding a List of Numbers

## 4.2.2. 递归的三法则

Like the robots of Asimov, all recursive algorithms must obey three important laws:

- A recursive algorithm must have a base case.
- A recursive algorithm must change its state and move toward the base case.
- A recursive algorithm must call itself, recursively.

Let's look at each one of these laws in more detail and see how it was used in the listsum algorithm. First, a base case is the condition that allows the algorithm to stop recursing. A base case is typically a problem that is small enough to solve directly. In the listsum algorithm the base case is a list of length 1.

To obey the second law, we must arrange for a change of state that moves the algorithm toward the base case. A change of state means that some data that the algorithm is using is modified. Usually the data that represents our problem gets smaller in some way. In the listsum algorithm our primary data structure is a list, so we must focus our state-changing efforts on the list. Since the base case is a list of length 1, a natural progression toward the base case is to shorten the list. This is exactly what happens on line 5 of ActiveCode 2 when we call listsum with a shorter list.

The final law is that the algorithm must call itself. This is the very definition of recursion. Recursion is a confusing concept to many beginning programmers. As a novice programmer, you have learned that functions are good because you can take a large problem and break it up into smaller problems. The smaller problems can be solved by writing a function to solve each problem. When we talk about recursion it may seem that we are talking ourselves in circles. We have a problem to solve with a function, but that function solves the problem by calling itself! But the logic is not circular at all; the logic of recursion is an elegant expression of solving a problem by breaking it down into a smaller and easier problems.

In the remainder of this chapter we will look at more examples of recursion. In each case we will focus on designing a solution to a problem by using the three laws of recursion.

### Self Check

Q-15: How many recursive calls are made when computing the sum of the list [2,4,6,8,10]?

- a) 6
- b) 5
- c) 4
- d) 3

Q-16: Suppose you are going to write a recursive function to calculate the factorial of a number. fact(n) returns  $n * n-1 * n-2 * \dots$ . Where the factorial of zero is defined to be 1. What would be the most appropriate base case?

- a)  $n == 0$
- b)  $n == 1$
- c)  $n >= 0$
- d)  $n <= 1$

### 4.2.3. 将整数转换为字符串形式的任意进制表示

Suppose you want to convert an integer to a string in some base between binary and hexadecimal. For example, convert the integer 10 to its string representation in decimal as "10", or to its string representation in binary as "1010". While there are many algorithms to solve this problem, including the algorithm discussed in the stack section, the recursive formulation of the problem is very elegant.

Let's look at a concrete example using base 10 and the number 769. Suppose we have a sequence of characters corresponding to the first 10 digits, like `convString = "0123456789"`. It is easy to convert a number less than 10 to its string equivalent by looking it up in the sequence. For example, if the number is 9, then the string is `convString[9]` or "9". If we can arrange to break up the number 769 into three single-digit numbers, 7, 6, and 9, then converting it to a string is simple. A number less than 10 sounds like a good base case.

Knowing what our base is suggests that the overall algorithm will involve three components:

- 1) Reduce the original number to a series of single-digit numbers.
- 2) Convert the single digit-number to a string using a lookup.
- 3) Concatenate the single-digit strings together to form the final result.

The next step is to figure out how to change state and make progress toward the base case. Since we are working with an integer, let's consider what mathematical operations might reduce a number. The most likely candidates are division and subtraction. While subtraction might work, it is unclear what we should subtract from what. Integer division with remainders gives us a clear direction. Let's look at what happens if we divide a number by the base we are trying to convert to.

Using integer division to divide 769 by 10, we get 76 with a remainder of 9. This gives us two good results. First, the remainder is a number less than our base that can be converted to a string immediately by lookup. Second, we get a number that is smaller than our original and moves us toward the base case of having a single number less than our base. Now our job is to convert 76 to its string representation. Again we will use integer division plus remainder to get results of 7 and 6 respectively. Finally, we have reduced the problem to converting 7, which we can do easily since it satisfies the base case condition of  $n < \text{base}$ , where  $\text{base} = 10$ . The series of operations we have just performed is illustrated in Figure 3. Notice that the numbers we want to remember are in the remainder boxes along the right side of the diagram.



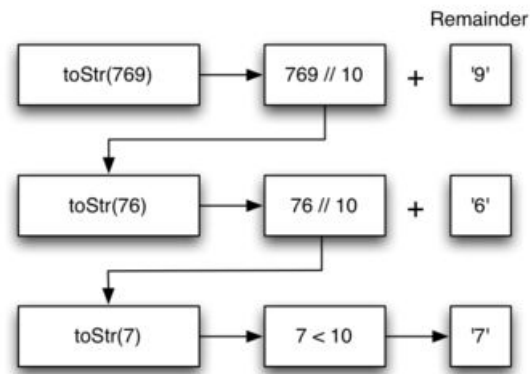


Figure 3: Converting an Integer to a String in Base 10

ActiveCode 1 shows the Python code that implements the algorithm outlined above for any base between 2 and 16.

```

def toStr(n,base):
    convertString = "0123456789ABCDEF"
    if n < base:
        return convertString[n]
    else:
        return toStr(n//base,base) + convertString[n%base]
print(toStr(1453,16))

```

#### Recursively Converting from Integer to String (lst\_rectostr)

Notice that in line 3 we check for the base case where  $n$  is less than the base we are converting to. When we detect the base case, we stop recursing and simply return the string from the `convertString` sequence. In line 6 we satisfy both the second and third laws – by making the recursive call and by reducing the problem size – using division.

Let's trace the algorithm again; this time we will convert the number 10 to its base 2 string representation ("1010").

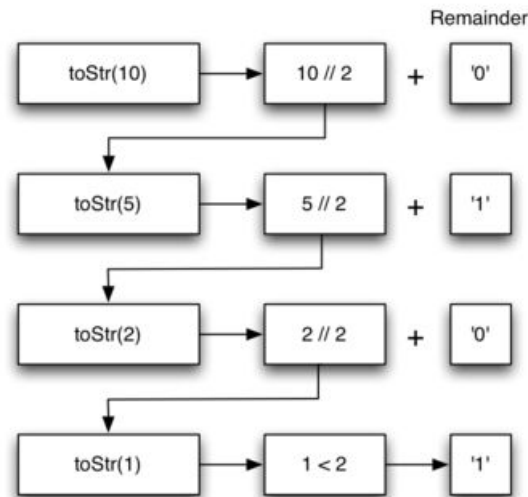


Figure 4: Converting the Number 10 to its Base 2 String Representation

Figure 4 shows that we get the results we are looking for, but it looks like the digits are in the wrong order. The algorithm works correctly because we make the recursive call first on line 6, then we add the string representation of the remainder. If we reversed returning the convertString lookup and returning the toStr call, the resulting string would be backward! But by delaying the concatenation operation until after the recursive call has returned, we get the result in the proper order. This should remind you of our discussion of stacks back in the previous chapter.

### Self Check

Write a function that takes a string as a parameter and returns a new string that is the reverse of the old string.

Write a function that takes a string as a parameter and returns True if the string is a palindrome, False otherwise. Remember that a string is a palindrome if it is spelled the same both forward and backward. For example: radar is a palindrome. for bonus points palindromes can also be phrases, but you need to remove the spaces and punctuation before checking. for example: madam i' m adam is a palindrome. Other fun palindromes include:

- kayak
- aibohphobia
- Live not on evil
- Reviled did I live, said I, as evil I did deliver
- Go hang a salami; I' m a lasagna hog.
- Able was I ere I saw Elba
- Kanakanak - a town in Alaska
- Wassamassaw - a town in South Dakota

## 4.3. 栈帧：实现递归

Suppose that instead of concatenating the result of the recursive call to toStr with the string from

convertString, we modified our algorithm to push the strings onto a stack prior to making the recursive call. The code for this modified algorithm is shown in ActiveCode 1.

```
from pythonds.basic.stack import Stack
rStack = Stack()
def toStr(n,base):
    convertString = "0123456789ABCDEF"
    while n > 0:
        if n < base:
            rStack.push(convertString[n])
        else:
            rStack.push(convertString[n % base])
        n = n // base
    res = ""
    while not rStack.isEmpty():
        res = res + str(rStack.pop())
    return res

print(toStr(1453,16))
```

#### Converting an Integer to a String Using a Stack (lst\_recstack)

Each time we make a call to toStr, we push a character on the stack. Returning to the previous example we can see that after the fourth call to toStr the stack would look like Figure 5. Notice that now we can simply pop the characters off the stack and concatenate them into the final result, "1010".

../\_images/recstack.png

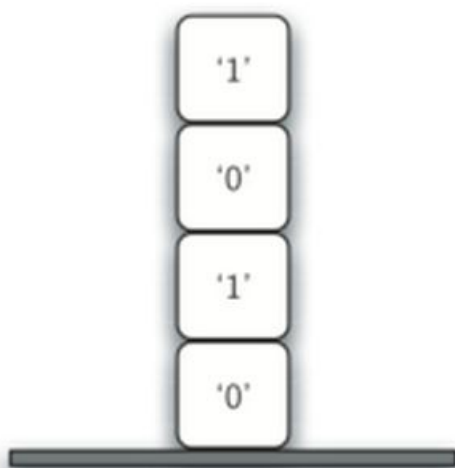


Figure 5: Strings Placed on the Stack During Conversion

The previous example gives us some insight into how Python implements a recursive function call. When a function is called in Python, a stack frame is allocated to handle the local variables of the

function. When the function returns, the return value is left on top of the stack for the calling function to access. Figure 6 illustrates the call stack after the return statement on line 4.

../\_images/newcallstack.png

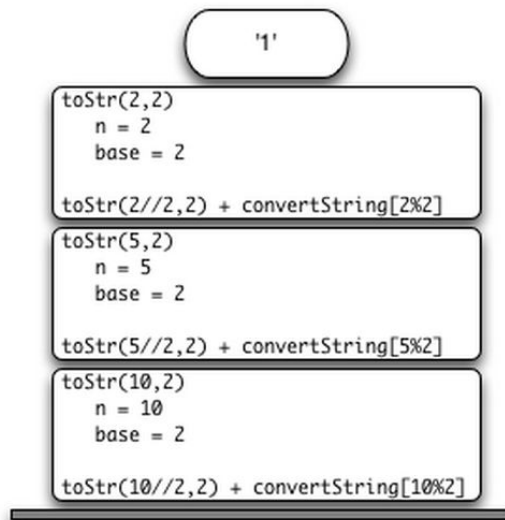


Figure 6: Call Stack Generated from toStr(10,2)

Notice that the call to toStr(2//2,2) leaves a return value of "1" on the stack. This return value is then used in place of the function call (toStr(1,2)) in the expression "1" + convertString[2%2], which will leave the string "10" on the top of the stack. In this way, the Python call stack takes the place of the stack we used explicitly in Listing 4. In our list summing example, you can think of the return value on the stack taking the place of an accumulator variable.

The stack frames also provide a scope for the variables used by the function. Even though we are calling the same function over and over, each call creates a new scope for the variables that are local to the function.

## 4.4. 图示递归

In the previous section we looked at some problems that were easy to solve using recursion; however, it can still be difficult to find a mental model or a way of visualizing what is happening in a recursive function. This can make recursion difficult for people to grasp. In this section we will look at a couple of examples of using recursion to draw some interesting pictures. As you watch these pictures take shape you will get some new insight into the recursive process that may be helpful in cementing your understanding of recursion.

The tool we will use for our illustrations is Python's turtle graphics module called turtle. The turtle module is standard with all versions of Python and is very easy to use. The metaphor is quite simple. You can create a turtle and the turtle can move forward, backward, turn left, turn right, etc. The turtle can have its tail up or down. When the turtle's tail is down and the turtle moves it draws a line as it moves. To increase the artistic value of the turtle you can change the width of the

tail as well as the color of the ink the tail is dipped in.

Here is a simple example to illustrate some turtle graphics basics. We will use the turtle module to draw a spiral recursively. ActiveCode 1 shows how it is done. After importing the turtle module we create a turtle. When the turtle is created it also creates a window for itself to draw in. Next we define the drawSpiral function. The base case for this simple function is when the length of the line we want to draw, as given by the len parameter, is reduced to zero or less. If the length of the line is longer than zero we instruct the turtle to go forward by len units and then turn right 90 degrees. The recursive step is when we call drawSpiral again with a reduced length. At the end of ActiveCode 1 you will notice that we call the function myWin.exitonclick(), this is a handy little method of the window that puts the turtle into a wait mode until you click inside the window, after which the program cleans up and exits.

```
import turtle

myTurtle = turtle.Turtle()
myWin = turtle.Screen()

def drawSpiral(myTurtle, lineLen):
    if lineLen > 0:
        myTurtle.forward(lineLen)
        myTurtle.right(90)
        drawSpiral(myTurtle, lineLen-5)

drawSpiral(myTurtle, 100)
myWin.exitonclick()
```

Drawing a Recursive Spiral using turtle (lst\_turt1)

That is really about all the turtle graphics you need to know in order to make some pretty impressive drawings. For our next program we are going to draw a fractal tree. Fractals come from a branch of mathematics, and have much in common with recursion. The definition of a fractal is that when you look at it the fractal has the same basic shape no matter how much you magnify it. Some examples from nature are the coastlines of continents, snowflakes, mountains, and even trees or shrubs. The fractal nature of many of these natural phenomenon makes it possible for programmers to generate very realistic looking scenery for computer generated movies. In our next example we will generate a fractal tree.

To understand how this is going to work it is helpful to think of how we might describe a tree using a fractal vocabulary. Remember that we said above that a fractal is something that looks the same at all different levels of magnification. If we translate this to trees and shrubs we might say that even a small twig has the same shape and characteristics as a whole tree. Using this idea we could say that a tree is a trunk, with a smaller tree going off to the right and another smaller tree going off to the left. If you think of this definition recursively it means that we will apply the

recursive definition of a tree to both of the smaller left and right trees.

Let's translate this idea to some Python code. Listing 1 shows how we can use our turtle to generate a fractal tree. Let's look at the code a bit more closely. You will see that on lines 5 and 7 we are making a recursive call. On line 5 we make the recursive call right after the turtle turns to the right by 20 degrees; this is the right tree mentioned above. Then in line 7 the turtle makes another recursive call, but this time after turning left by 40 degrees. The reason the turtle must turn left by 40 degrees is that it needs to undo the original 20 degree turn to the right and then do an additional 20 degree turn to the left in order to draw the left tree. Also notice that each time we make a recursive call to tree we subtract some amount from the branchLen parameter; this is to make sure that the recursive trees get smaller and smaller. You should also recognize the initial if statement on line 2 as a check for the base case of branchLen getting too small.

Listing 1

```
def tree(branchLen,t):
    if branchLen > 5:
        t.forward(branchLen)
        t.right(20)
        tree(branchLen-15,t)
        t.left(40)
        tree(branchLen-10,t)
        t.right(20)
        t.backward(branchLen)
```

The complete program for this tree example is shown in ActiveCode 2. Before you run the code think about how you expect to see the tree take shape. Look at the recursive calls and think about how this tree will unfold. Will it be drawn symmetrically with the right and left halves of the tree taking shape simultaneously? Will it be drawn right side first then left side?

```
import turtle

def tree(branchLen,t):
    if branchLen > 5:
        t.forward(branchLen)
        t.right(20)
        tree(branchLen-15,t)
        t.left(40)
        tree(branchLen-15,t)
        t.right(20)
        t.backward(branchLen)

def main():
    t = turtle.Turtle()
```

```
myWin = turtle.Screen()
t.left(90)
t.up()
t.backward(100)
t.down()
t.color("green")
tree(75,t)
myWin.exitonclick()

main()
```

#### Recursively Drawing a Tree (lst\_complete\_tree)

Notice how each branch point on the tree corresponds to a recursive call, and notice how the tree is drawn to the right all the way down to its shortest twig. You can see this in Figure 1. Now, notice how the program works its way back up the trunk until the entire right side of the tree is drawn. You can see the right half of the tree in Figure 2. Then the left side of the tree is drawn, but not by going as far out to the left as possible. Rather, once again the entire right side of the left tree is drawn until we finally make our way out to the smallest twig on the left.

../\_images/tree1.png

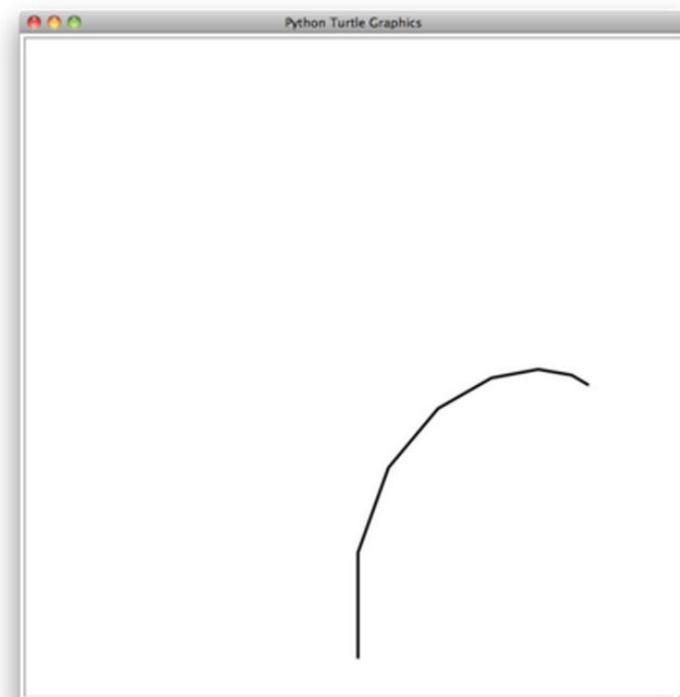


Figure 1: The Beginning of a Fractal Tree

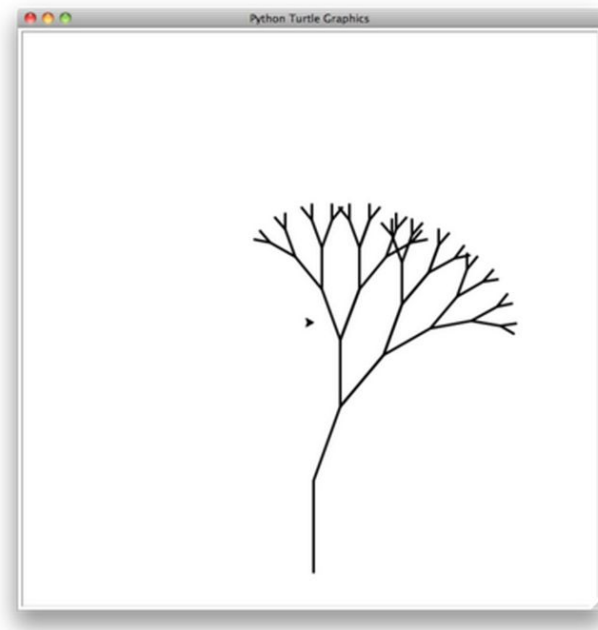


Figure 2: The First Half of the Tree

This simple tree program is just a starting point for you, and you will notice that the tree does not look particularly realistic because nature is just not as symmetric as a computer program. The exercises at the end of the chapter will give you some ideas for how to explore some interesting options to make your tree look more realistic.

### Self Check

Modify the recursive tree program using one or all of the following ideas:

- Modify the thickness of the branches so that as the branchLen gets smaller, the line gets thinner.
- Modify the color of the branches so that as the branchLen gets very short it is colored like a leaf.
- Modify the angle used in turning the turtle so that at each branch point the angle is selected at random in some range. For example choose the angle between 15 and 45 degrees. Play around to see what looks good.
- Modify the branchLen recursively so that instead of always subtracting the same amount you subtract a random amount in some range.

### 4.4.1. 谢尔宾斯基三角形 Sierpinski Triangle

Another fractal that exhibits the property of self-similarity is the Sierpinski triangle. An example is shown in Figure 3. The Sierpinski triangle illustrates a three-way recursive algorithm. The procedure for drawing a Sierpinski triangle by hand is simple. Start with a single large triangle. Divide this large triangle into four new triangles by connecting the midpoint of each side. Ignoring the middle triangle that you just created, apply the same procedure to each of the three corner triangles. Each time you create a new set of triangles, you recursively apply this procedure to the three smaller corner triangles. You can continue to apply this procedure indefinitely if you have a



sharp enough pencil. Before you continue reading, you may want to try drawing the Sierpinski triangle yourself, using the method described.

../\_images/sierpinski.png

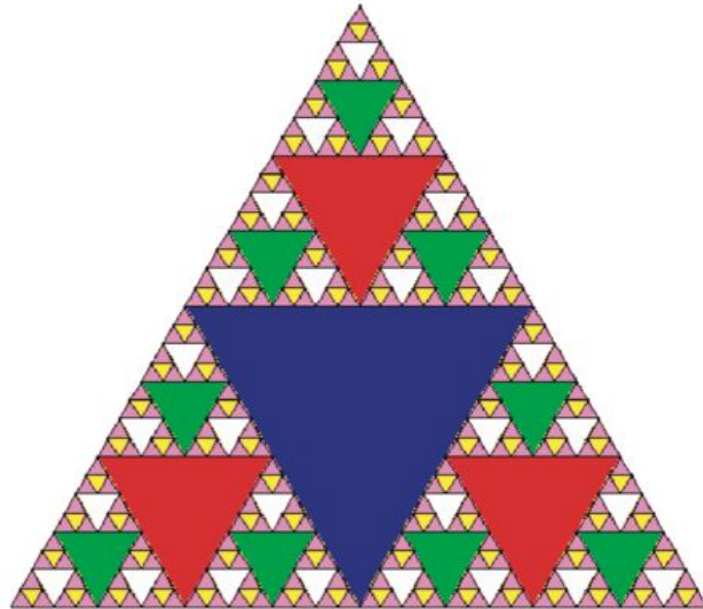


Figure 3: The Sierpinski Triangle

Since we can continue to apply the algorithm indefinitely, what is the base case? We will see that the base case is set arbitrarily as the number of times we want to divide the triangle into pieces. Sometimes we call this number the “degree” of the fractal. Each time we make a recursive call, we subtract 1 from the degree until we reach 0. When we reach a degree of 0, we stop making recursive calls. The code that generated the Sierpinski Triangle in Figure 3 is shown in ActiveCode 1.

```
import turtle

def drawTriangle(points,color,myTurtle):
    myTurtle.fillcolor(color)
    myTurtle.up()
    myTurtle.goto(points[0][0],points[0][1])
    myTurtle.down()
    myTurtle.begin_fill()
    myTurtle.goto(points[1][0],points[1][1])
    myTurtle.goto(points[2][0],points[2][1])
    myTurtle.goto(points[0][0],points[0][1])
    myTurtle.end_fill()

def getMid(p1,p2):
    return ( (p1[0]+p2[0]) / 2, (p1[1] + p2[1]) / 2)
```

```

def sierpinski(points,degree,myTurtle):
    colormap = ['blue','red','green','white','yellow', \
                'violet','orange']
    drawTriangle(points,colormap[degree],myTurtle)
    if degree > 0:
        sierpinski([points[0], \
                    getMid(points[0], points[1]), \
                    getMid(points[0], points[2])], \
                    degree-1, myTurtle)
        sierpinski([points[1], \
                    getMid(points[0], points[1]), \
                    getMid(points[1], points[2])], \
                    degree-1, myTurtle)
        sierpinski([points[2], \
                    getMid(points[2], points[1]), \
                    getMid(points[0], points[2])], \
                    degree-1, myTurtle)

def main():
    myTurtle = turtle.Turtle()
    myWin = turtle.Screen()
    myPoints = [[-100,-50],[0,100],[100,-50]]
    sierpinski(myPoints,3,myTurtle)
    myWin.exitonclick()

main()

```

#### Drawing a Sierpinski Triangle (1st\_st)

The program in ActiveCode 1 follows the ideas outlined above. The first thing `sierpinski` does is draw the outer triangle. Next, there are three recursive calls, one for each of the new corner triangles we get when we connect the midpoints. Once again we make use of the standard `turtle` module that comes with Python. You can learn all the details of the methods available in the `turtle` module by using `help('turtle')` from the Python prompt.

Look at the code and think about the order in which the triangles will be drawn. While the exact order of the corners depends upon how the initial set is specified, let's assume that the corners are ordered lower left, top, lower right. Because of the way the `sierpinski` function calls itself, `sierpinski` works its way to the smallest allowed triangle in the lower-left corner, and then begins to fill out the rest of the triangles working back. Then it fills in the triangles in the top corner by working toward the smallest, topmost triangle. Finally, it fills in the lower-right corner, working its way toward the smallest triangle in the lower right.

Sometimes it is helpful to think of a recursive algorithm in terms of a diagram of function calls.

Figure 4 shows that the recursive calls are always made going to the left. The active functions are outlined in black, and the inactive function calls are in gray. The farther you go toward the bottom of Figure 4, the smaller the triangles. The function finishes drawing one level at a time; once it is finished with the bottom left it moves to the bottom middle, and so on.

../\_images/stCallTree.png

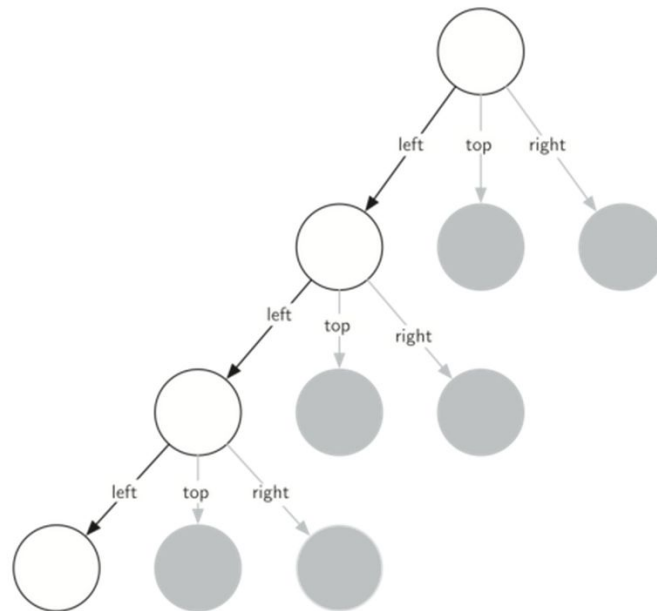


Figure 4: Building a Sierpinski Triangle

The `sierpinski` function relies heavily on the `getMid` function. `getMid` takes as arguments two endpoints and returns the point halfway between them. In addition, ActiveCode 1 has a function that draws a filled triangle using the `begin_fill` and `end_fill` turtle methods.

## 4.5. 复杂递归问题

In the previous sections we looked at some problems that are relatively easy to solve and some graphically interesting problems that can help us gain a mental model of what is happening in a recursive algorithm. In this section we will look at some problems that are really difficult to solve using an iterative programming style but are very elegant and easy to solve using recursion. We will finish up by looking at a deceptive problem that at first looks like it has an elegant recursive solution but in fact does not.

### 4.5.1. 河内塔问题 Towers of Hanoi

The Tower of Hanoi puzzle was invented by the French mathematician Edouard Lucas in 1883. He was inspired by a legend that tells of a Hindu temple where the puzzle was presented to young priests. At the beginning of time, the priests were given three poles and a stack of 64 gold disks, each disk a little smaller than the one beneath it. Their assignment was to transfer all 64 disks

from one of the three poles to another, with two important constraints. They could only move one disk at a time, and they could never place a larger disk on top of a smaller one. The priests worked very efficiently, day and night, moving one disk every second. When they finished their work, the legend said, the temple would crumble into dust and the world would vanish.

Although the legend is interesting, you need not worry about the world ending any time soon. The number of moves required to correctly move a tower of 64 disks is  $2^{64} - 1 = 18,446,744,073,709,551,615$ . At a rate of one move per second, that is 584,942,417,355 years! Clearly there is more to this puzzle than meets the eye.

Figure 1 shows an example of a configuration of disks in the middle of a move from the first peg to the third. Notice that, as the rules specify, the disks on each peg are stacked so that smaller disks are always on top of the larger disks. If you have not tried to solve this puzzle before, you should try it now. You do not need fancy disks and poles – a pile of books or pieces of paper will work.

Image

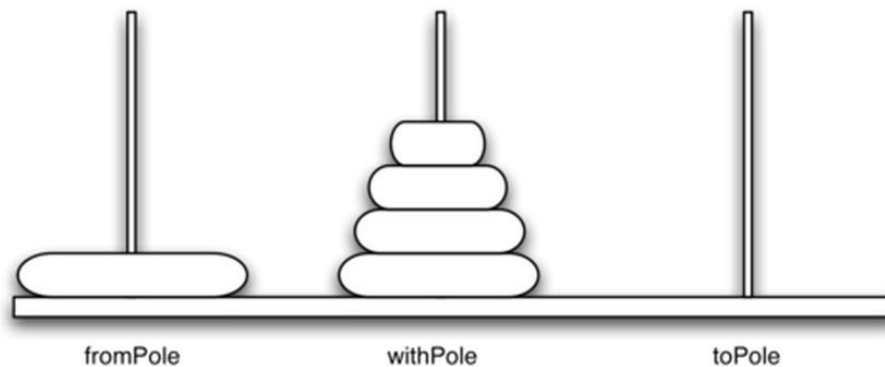


Figure 1: An Example Arrangement of Disks for the Tower of Hanoi

How do we go about solving this problem recursively? How would you go about solving this problem at all? What is our base case? Let's think about this problem from the bottom up. Suppose you have a tower of five disks, originally on peg one. If you already knew how to move a tower of four disks to peg two, you could then easily move the bottom disk to peg three, and then move the tower of four from peg two to peg three. But what if you do not know how to move a tower of height four? Suppose that you knew how to move a tower of height three to peg three; then it would be easy to move the fourth disk to peg two and move the three from peg three on top of it. But what if you do not know how to move a tower of three? How about moving a tower of two disks to peg two and then moving the third disk to peg three, and then moving the tower of height two on top of it? But what if you still do not know how to do this? Surely you would agree that moving a single disk to peg three is easy enough, trivial you might even say. This sounds like a base case in the making.

Here is a high-level outline of how to move a tower from the starting pole, to the goal pole, using an intermediate pole:

- Move a tower of height-1 to an intermediate pole, using the final pole.
- Move the remaining disk to the final pole.
- Move the tower of height-1 from the intermediate pole to the final pole using the original pole.

As long as we always obey the rule that the larger disks remain on the bottom of the stack, we can use the three steps above recursively, treating any larger disks as though they were not even there. The only thing missing from the outline above is the identification of a base case. The simplest Tower of Hanoi problem is a tower of one disk. In this case, we need move only a single disk to its final destination. A tower of one disk will be our base case. In addition, the steps outlined above move us toward the base case by reducing the height of the tower in steps 1 and 3. Listing 1 shows the Python code to solve the Tower of Hanoi puzzle.

Listing 1

```
def moveTower(height,fromPole, toPole, withPole):
    if height >= 1:
        moveTower(height-1,fromPole,withPole,toPole)
        moveDisk(fromPole,toPole)
        moveTower(height-1,withPole,toPole,fromPole)
```

Notice that the code in Listing 1 is almost identical to the English description. The key to the simplicity of the algorithm is that we make two different recursive calls, one on line 3 and a second on line 5. On line 3 we move all but the bottom disk on the initial tower to an intermediate pole. The next line simply moves the bottom disk to its final resting place. Then on line 5 we move the tower from the intermediate pole to the top of the largest disk. The base case is detected when the tower height is 0; in this case there is nothing to do, so the moveTower function simply returns. The important thing to remember about handling the base case this way is that simply returning from moveTower is what finally allows the moveDisk function to be called.

The function moveDisk, shown in Listing 2, is very simple. All it does is print out that it is moving a disk from one pole to another. If you type in and run the moveTower program you can see that it gives you a very efficient solution to the puzzle.

Listing 2

```
def moveDisk(fp,tp):
    print("moving disk from",fp,"to",tp)
```

The program in ActiveCode 1 provides the entire solution for three disks.

```
def moveTower(height,fromPole, toPole, withPole):
    if height >= 1:
        moveTower(height-1,fromPole,withPole,toPole)
        moveDisk(fromPole,toPole)
```

```
        moveTower(height-1,withPole,toPole,fromPole)

def moveDisk(fp,tp):
    print("moving disk from",fp,"to",tp)

moveTower(3,"A","B","C")
```

Solving Tower of Hanoi Recursively (hanoi)

Now that you have seen the code for both `moveTower` and `moveDisk`, you may be wondering why we do not have a data structure that explicitly keeps track of what disks are on what poles. Here is a hint: if you were going to explicitly keep track of the disks, you would probably use three `Stack` objects, one for each pole. The answer is that Python provides the stacks that we need implicitly through the call stack.

## 4.6. 探索迷宫

In this section we will look at a problem that has relevance to the expanding world of robotics: How do you find your way out of a maze? If you have a Roomba vacuum cleaner for your dorm room (don't all college students?) you will wish that you could reprogram it using what you have learned in this section. The problem we want to solve is to help our turtle find its way out of a virtual maze. The maze problem has roots as deep as the Greek myth about Theseus who was sent into a maze to kill the minotaur. Theseus used a ball of thread to help him find his way back out again once he had finished off the beast. In our problem we will assume that our turtle is dropped down somewhere into the middle of the maze and must find its way out. Look at Figure 2 to get an idea of where we are going in this section.

../\_images/maze.png

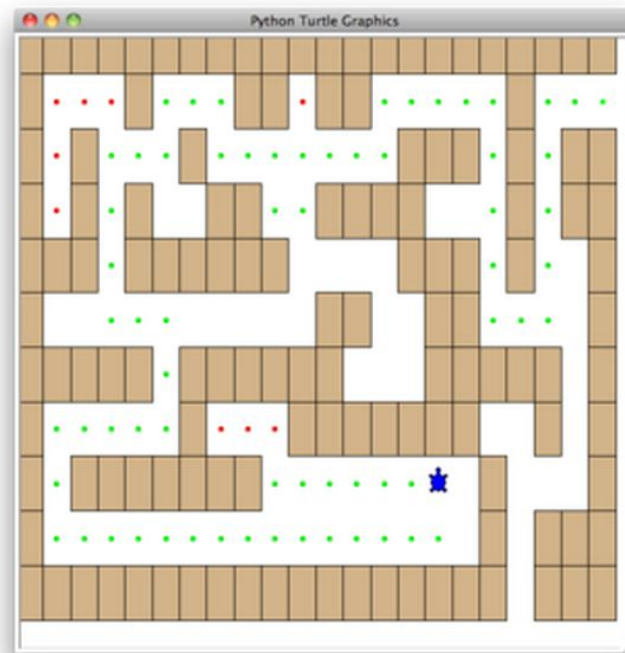


Figure 2: The Finished Maze Search Program

To make it easier for us we will assume that our maze is divided up into “squares.” Each square of the maze is either open or occupied by a section of wall. The turtle can only pass through the open squares of the maze. If the turtle bumps into a wall it must try a different direction. The turtle will require a systematic procedure to find its way out of the maze. Here is the procedure:

- From our starting position we will first try going North one square and then recursively try our procedure from there.
- If we are not successful by trying a Northern path as the first step then we will take a step to the South and recursively repeat our procedure.
- If South does not work then we will try a step to the West as our first step and recursively apply our procedure.
- If North, South, and West have not been successful then apply the procedure recursively from a position one step to our East.
- If none of these directions works then there is no way to get out of the maze and we fail.

Now, that sounds pretty easy, but there are a couple of details to talk about first. Suppose we take our first recursive step by going North. By following our procedure our next step would also be to the North. But if the North is blocked by a wall we must look at the next step of the procedure and try going to the South. Unfortunately that step to the south brings us right back to our original starting place. If we apply the recursive procedure from there we will just go back one step to the North and be in an infinite loop. So, we must have a strategy to remember where we have been. In this case we will assume that we have a bag of bread crumbs we can drop along our way. If we take a step in a certain direction and find that there is a bread crumb already on that square, we know that we should immediately back up and try the next direction in our procedure. As we will see when we look at the code for this algorithm, backing up is as simple as returning from a recursive function call.

As we do for all recursive algorithms let us review the base cases. Some of them you may already have guessed based on the description in the previous paragraph. In this algorithm, there are four base cases to consider:

- The turtle has run into a wall. Since the square is occupied by a wall no further exploration can take place.
- The turtle has found a square that has already been explored. We do not want to continue exploring from this position or we will get into a loop.
- We have found an outside edge, not occupied by a wall. In other words we have found an exit from the maze.
- We have explored a square unsuccessfully in all four directions.

For our program to work we will need to have a way to represent the maze. To make this even more interesting we are going to use the turtle module to draw and explore our maze so we can watch this algorithm in action. The maze object will provide the following methods for us to use in writing our search algorithm:

- `__init__` Reads in a data file representing a maze, initializes the internal representation of the maze, and finds the starting position for the turtle.
- `drawMaze` Draws the maze in a window on the screen.
- `updatePosition` Updates the internal representation of the maze and changes the position of the turtle in the window.
- `isExit` Checks to see if the current position is an exit from the maze.

The Maze class also overloads the index operator `[]` so that our algorithm can easily access the status of any particular square.

Let's examine the code for the search function which we call `searchFrom`. The code is shown in Listing 3. Notice that this function takes three parameters: a maze object, the starting row, and the starting column. This is important because as a recursive function the search logically starts again with each recursive call.

Listing 3

```
def searchFrom(maze, startRow, startColumn):
    maze.updatePosition(startRow, startColumn)
    # Check for base cases:
    # 1. We have run into an obstacle, return false
    if maze[startRow][startColumn] == OBSTACLE :
        return False
    # 2. We have found a square that has already been explored
    if maze[startRow][startColumn] == TRIED:
        return False
    # 3. Success, an outside edge not occupied by an obstacle
    if maze.isExit(startRow,startColumn):
```



```

        maze.updatePosition(startRow, startColumn, PART_OF_PATH)
        return True
    maze.updatePosition(startRow, startColumn, TRIED)

    # Otherwise, use logical short circuiting to try each
    # direction in turn (if needed)
    found = searchFrom(maze, startRow-1, startColumn) or \
            searchFrom(maze, startRow+1, startColumn) or \
            searchFrom(maze, startRow, startColumn-1) or \
            searchFrom(maze, startRow, startColumn+1)
    if found:
        maze.updatePosition(startRow, startColumn, PART_OF_PATH)
    else:
        maze.updatePosition(startRow, startColumn, DEAD_END)
    return found

```

As you look through the algorithm you will see that the first thing the code does (line 2) is call `updatePosition`. This is simply to help you visualize the algorithm so that you can watch exactly how the turtle explores its way through the maze. Next the algorithm checks for the first three of the four base cases: Has the turtle run into a wall (line 5)? Has the turtle circled back to a square already explored (line 8)? Has the turtle found an exit (line 11)? If none of these conditions is true then we continue the search recursively.

You will notice that in the recursive step there are four recursive calls to `searchFrom`. It is hard to predict how many of these recursive calls will be used since they are all connected by `or` statements. If the first call to `searchFrom` returns `True` then none of the last three calls would be needed. You can interpret this as meaning that a step to (row-1,column) (or North if you want to think geographically) is on the path leading out of the maze. If there is not a good path leading out of the maze to the North then the next recursive call is tried, this one to the South. If South fails then try West, and finally East. If all four recursive calls return false then we have found a dead end. You should download or type in the whole program and experiment with it by changing the order of these calls.

The code for the `Maze` class is shown in Listing 4, Listing 5, and Listing 6. The `__init__` method takes the name of a file as its only parameter. This file is a text file that represents a maze by using “+” characters for walls, spaces for open squares, and the letter “S” to indicate the starting position. Figure 3 is an example of a maze data file. The internal representation of the maze is a list of lists. Each row of the `mazelist` instance variable is also a list. This secondary list contains one character per square using the characters described above. For the data file in Figure 3 the internal representation looks like the following:

```

[ ['+', '+', '+', '+', ..., '+', '+', '+', '+', '+', '+', '+'],
  ['+', ' ', ' ', ' ', ' ', ..., ' ', ' ', ' ', '+', ' ', ' ', ' '],
  ['+', ' ', '+', ' ', ' ', ..., '+', '+', ' ', '+', ' ', '+', '+'],

```

```
[ '+', ' ', '+', ' ', ..., ' ', ' ', ' ', '+', ' ', '+', '+' ],
[ '+', '+', '+', ' ', ..., '+', '+', ' ', '+', ' ', ' ', '+' ],
[ '+', ' ', ' ', ' ', ..., '+', '+', ' ', ' ', ' ', ' ', '+' ],
[ '+', '+', '+', '+', ..., '+', '+', '+', '+', '+', ' ', '+' ],
[ '+', ' ', ' ', ' ', ..., '+', '+', ' ', ' ', ' ', '+', '+' ],
[ '+', ' ', '+', '+', ..., ' ', ' ', '+', ' ', ' ', ' ', '+' ],
[ '+', ' ', ' ', ' ', ..., ' ', ' ', '+', ' ', '+', '+', '+' ],
[ '+', '+', '+', '+', ..., '+', '+', '+', ' ', '+', '+', '+' ]]
```

The drawMaze method uses this internal representation to draw the initial view of the maze on the screen.

Figure 3: An Example Maze Data File

```
+++++
+  +  ++ ++  +
+ +  +      +++ + ++
+ + + ++ +++++ + ++
+++ ++++++  +++ + +
+          ++ ++  +
+++++ ++++++ +++++ +
+      +  ++++++ + +
+ ++++++      S +  +
+          + +++
+++++ ++++++ +++
```

The updatePosition method, as shown in Listing 5 uses the same internal representation to see if the turtle has run into a wall. It also updates the internal representation with a “.” or “-” to indicate that the turtle has visited a particular square or if the square is part of a dead end. In addition, the updatePosition method uses two helper methods, moveTurtle and dropBreadCrumb, to update the view on the screen.

Finally, the isExit method uses the current position of the turtle to test for an exit condition. An exit condition is whenever the turtle has navigated to the edge of the maze, either row zero or column zero, or the far right column or the bottom row.

Listing 4

```
class Maze:
    def __init__(self,mazeFileName):
        rowsInMaze = 0
        columnsInMaze = 0
        self.mazelist = []
        mazeFile = open(mazeFileName,'r')
        rowsInMaze = 0
        for line in mazeFile:
            rowList = []
```

```

        col = 0
        for ch in line[:-1]:
            rowList.append(ch)
            if ch == 'S':
                self.startRow = rowsInMaze
                self.startCol = col
            col = col + 1
        rowsInMaze = rowsInMaze + 1
        self.mazelist.append(rowList)
        columnsInMaze = len(rowList)

    self.rowsInMaze = rowsInMaze
    self.columnsInMaze = columnsInMaze
    self.xTranslate = -columnsInMaze/2
    self.yTranslate = rowsInMaze/2
    self.t = Turtle(shape='turtle')
    setup(width=600,height=600)
    setworldcoordinates(-(columnsInMaze-1)/2-.5,
                        -(rowsInMaze-1)/2-.5,
                        (columnsInMaze-1)/2+.5,
                        (rowsInMaze-1)/2+.5)

```

Listing 5

```

def drawMaze(self):
    for y in range(self.rowsInMaze):
        for x in range(self.columnsInMaze):
            if self.mazelist[y][x] == OBSTACLE:
                self.drawCenteredBox(x+self.xTranslate,
                                    -y+self.yTranslate,
                                    'tan')

    self.t.color('black','blue')

def drawCenteredBox(self,x,y,color):
    tracer(0)
    self.t.up()
    self.t.goto(x-.5,y-.5)
    self.t.color('black',color)
    self.t.setheading(90)
    self.t.down()
    self.t.begin_fill()
    for i in range(4):
        self.t.forward(1)
        self.t.right(90)

```

```

        self.t.end_fill()
        update()
        tracer(1)

    def moveTurtle(self,x,y):
        self.t.up()
        self.t.setheading(self.t.towards(x+self.xTranslate,
                                          -y+self.yTranslate))
        self.t.goto(x+self.xTranslate,-y+self.yTranslate)

    def dropBreadcrumb(self,color):
        self.t.dot(color)

    def updatePosition(self,row,col,val=None):
        if val:
            self.mazelist[row][col] = val
            self.moveTurtle(col,row)

            if val == PART_OF_PATH:
                color = 'green'
            elif val == OBSTACLE:
                color = 'red'
            elif val == TRIED:
                color = 'black'
            elif val == DEAD_END:
                color = 'red'
            else:
                color = None

            if color:
                self.dropBreadcrumb(color)

```

Listing 6

```

def isExit(self,row,col):
    return (row == 0 or
            row == self.rowsInMaze-1 or
            col == 0 or
            col == self.columnsInMaze-1 )

def __getitem__(self,idx):
    return self.mazelist[idx]

```

The complete program is shown in ActiveCode 1. This program uses the data file maze2.txt shown below. Note that it is a much more simple example file in that the exit is very close to the starting

position of the turtle.

```
+++++
+  +  ++ ++      +
      +      +++++
+ +  ++  ++++ ++ ++
+ +  + + ++  +++ +
+           ++ ++ ++
+++++ + +      ++ ++
+++++ +++ + + ++ +
+           + + S+ + +
+++++ +  + + +  + +
+++++
```

```
import turtle

PART_OF_PATH = 'O'
TRIED = '.'
OBSTACLE = '+'
DEAD_END = '-'

class Maze:
    def __init__(self,mazeFileName):
        rowsInMaze = 0
        columnsInMaze = 0
        self.mazelist = []
        mazeFile = open(mazeFileName,'r')
        rowsInMaze = 0
        for line in mazeFile:
            rowList = []
            col = 0
            for ch in line[:-1]:
                rowList.append(ch)
                if ch == 'S':
                    self.startRow = rowsInMaze
                    self.startCol = col
                col = col + 1
            rowsInMaze = rowsInMaze + 1
            self.mazelist.append(rowList)
            columnsInMaze = len(rowList)

        self.rowsInMaze = rowsInMaze
        self.columnsInMaze = columnsInMaze
        self.xTranslate = -columnsInMaze/2
```

```
self.yTranslate = rowsInMaze/2
self.t = turtle.Turtle()
self.t.shape('turtle')
```

Complete Maze Solver (completemaze)

### Self Check

Modify the maze search program so that the calls to searchFrom are in a different order. Watch the program run. Can you explain why the behavior is different? Can you predict what path the turtle will follow for a given change in order?

## 4.7. 动态规划

Many programs in computer science are written to optimize some value; for example, find the shortest path between two points, find the line that best fits a set of points, or find the smallest set of objects that satisfies some criteria. There are many strategies that computer scientists use to solve these problems. One of the goals of this book is to expose you to several different problem solving strategies. Dynamic programming is one strategy for these types of optimization problems.

A classic example of an optimization problem involves making change using the fewest coins. Suppose you are a programmer for a vending machine manufacturer. Your company wants to streamline effort by giving out the fewest possible coins in change for each transaction. Suppose a customer puts in a dollar bill and purchases an item for 37 cents. What is the smallest number of coins you can use to make change? The answer is six coins: two quarters, one dime, and three pennies. How did we arrive at the answer of six coins? We start with the largest coin in our arsenal (a quarter) and use as many of those as possible, then we go to the next lowest coin value and use as many of those as possible. This first approach is called a greedy method because we try to solve as big a piece of the problem as possible right away.

The greedy method works fine when we are using U.S. coins, but suppose that your company decides to deploy its vending machines in Lower Elbonia where, in addition to the usual 1, 5, 10, and 25 cent coins they also have a 21 cent coin. In this instance our greedy method fails to find the optimal solution for 63 cents in change. With the addition of the 21 cent coin the greedy method would still find the solution to be six coins. However, the optimal answer is three 21 cent pieces.

Let's look at a method where we could be sure that we would find the optimal answer to the problem. Since this section is about recursion, you may have guessed that we will use a recursive solution. Let's start with identifying the base case. If we are trying to make change for the same amount as the value of one of our coins, the answer is easy, one coin.

If the amount does not match we have several options. What we want is the minimum of a penny plus the number of coins needed to make change for the original amount minus a penny, or a nickel plus the number of coins needed to make change for the original amount minus five cents, or a dime plus the number of coins needed to make change for the original amount minus ten cents,

and so on. So the number of coins needed to make change for the original amount can be computed according to the following:

$$\text{numCoins} = \min \begin{cases} 1 + \text{numCoins}(\text{originalamount} - 1) \\ 1 + \text{numCoins}(\text{originalamount} - 5) \\ 1 + \text{numCoins}(\text{originalamount} - 10) \\ 1 + \text{numCoins}(\text{originalamount} - 25) \end{cases}$$

The algorithm for doing what we have just described is shown in Listing 7. In line 3 we are checking our base case; that is, we are trying to make change in the exact amount of one of our coins. If we do not have a coin equal to the amount of change, we make recursive calls for each different coin value less than the amount of change we are trying to make. Line 6 shows how we filter the list of coins to those less than the current value of change using a list comprehension. The recursive call also reduces the total amount of change we need to make by the value of the coin selected. The recursive call is made in line 7. Notice that on that same line we add 1 to our number of coins to account for the fact that we are using a coin. Just adding 1 is the same as if we had made a recursive call asking where we satisfy the base case condition immediately.

Listing 7

```
def recMC(coinValueList,change):
    minCoins = change
    if change in coinValueList:
        return 1
    else:
        for i in [c for c in coinValueList if c <= change]:
            numCoins = 1 + recMC(coinValueList,change-i)
            if numCoins < minCoins:
                minCoins = numCoins
        return minCoins

print(recMC([1,5,10,25],63))
```

The trouble with the algorithm in Listing 7 is that it is extremely inefficient. In fact, it takes 67,716,925 recursive calls to find the optimal solution to the 4 coins, 63 cents problem! To understand the fatal flaw in our approach look at Figure 5, which illustrates a small fraction of the 377 function calls needed to find the optimal set of coins to make change for 26 cents.

Each node in the graph corresponds to a call to `recMC`. The label on the node indicates the amount of change for which we are computing the number of coins. The label on the arrow indicates the coin that we just used. By following the graph we can see the combination of coins that got us to any point in the graph. The main problem is that we are re-doing too many calculations. For example, the graph shows that the algorithm would recalculate the optimal number of coins to make change for 15 cents at least three times. Each of these computations to find the optimal number of coins for 15 cents itself takes 52 function calls. Clearly we are wasting a lot of time and

effort recalculating old results.

Image

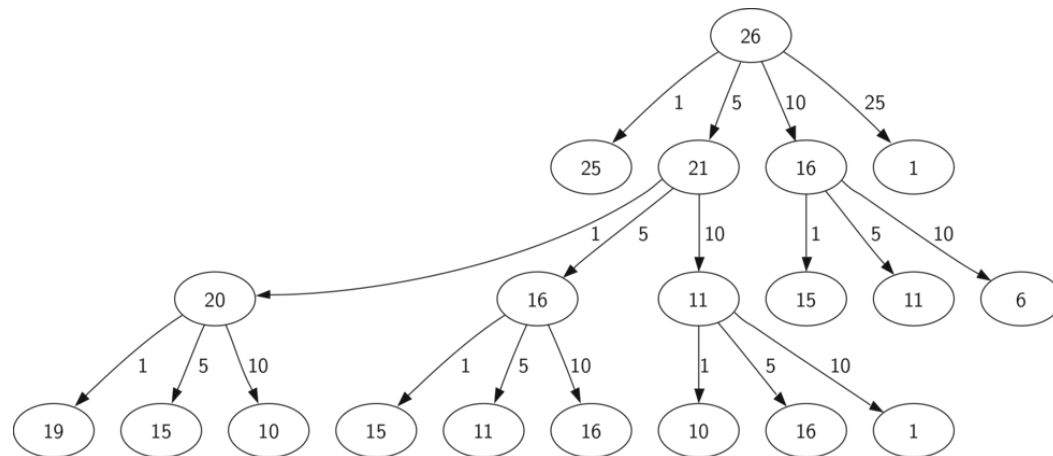


Figure 3: Call Tree for Listing 7

The key to cutting down on the amount of work we do is to remember some of the past results so we can avoid recomputing results we already know. A simple solution is to store the results for the minimum number of coins in a table when we find them. Then before we compute a new minimum, we first check the table to see if a result is already known. If there is already a result in the table, we use the value from the table rather than recomputing. ActiveCode 1 shows a modified algorithm to incorporate our table lookup scheme.

Run Save Load

```
def recDC(coinValueList,change,knownResults):
    minCoins = change
    if change in coinValueList:
        knownResults[change] = 1
        return 1
    elif knownResults[change] > 0:
        return knownResults[change]
    else:
        for i in [c for c in coinValueList if c <= change]:
            numCoins = 1 + recDC(coinValueList, change-i,
                                knownResults)
            if numCoins < minCoins:
                minCoins = numCoins
                knownResults[change] = minCoins
        return minCoins

print(recDC([1,5,10,25],63,[0]*64))
```



### Recursively Counting Coins with Table Lookup (lst\_change2)

Notice that in line 6 we have added a test to see if our table contains the minimum number of coins for a certain amount of change. If it does not, we compute the minimum recursively and store the computed minimum in the table. Using this modified algorithm reduces the number of recursive calls we need to make for the four coin, 63 cent problem to 221 calls!

Although the algorithm in AcitveCode 3 is correct, it looks and feels like a bit of a hack. Also, if we look at the knownResults lists we can see that there are some holes in the table. In fact the term for what we have done is not dynamic programming but rather we have improved the performance of our program by using a technique known as “memoization,” or more commonly called “caching.”

A truly dynamic programming algorithm will take a more systematic approach to the problem. Our dynamic programming solution is going to start with making change for one cent and systematically work its way up to the amount of change we require. This guarantees us that at each step of the algorithm we already know the minimum number of coins needed to make change for any smaller amount.

Let’s look at how we would fill in a table of minimum coins to use in making change for 11 cents. Figure 4 illustrates the process. We start with one cent. The only solution possible is one coin (a penny). The next row shows the minimum for one cent and two cents. Again, the only solution is two pennies. The fifth row is where things get interesting. Now we have two options to consider, five pennies or one nickel. How do we decide which is best? We consult the table and see that the number of coins needed to make change for four cents is four, plus one more penny to make five, equals five coins. Or we can look at zero cents plus one more nickel to make five cents equals 1 coin. Since the minimum of one and five is one we store 1 in the table. Fast forward again to the end of the table and consider 11 cents. Figure 5 shows the three options that we have to consider:

1. A penny plus the minimum number of coins to make change for  $11-1=10$  cents (1)
2. A nickel plus the minimum number of coins to make change for  $11-5=6$  cents (2)
3. A dime plus the minimum number of coins to make change for  $11-10=1$  cent (1)

Either option 1 or 3 will give us a total of two coins which is the minimum number of coins for 11 cents.

Image

Change to Make		1	2	3	4	5	6	6	8	9	10	11
Step of the Algorithm	1											
	1	2										
	1	2	3									
	1	2	3	4								
	1	2	3	4	1							
	...											
	1	2	3	4	1	2	3	4	5	1		
	1	2	3	4	1	2	3	4	5	1	2	

Figure 4: Minimum Number of Coins Needed to Make Change

Image

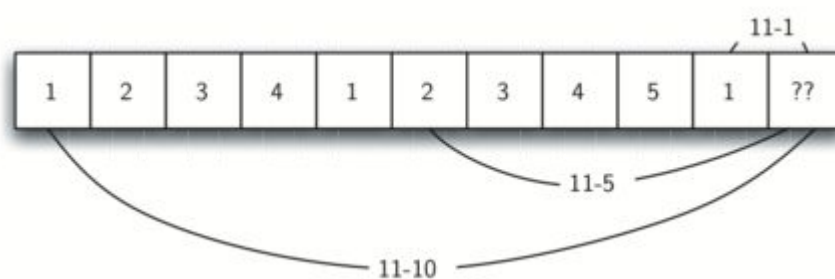


Figure 5: Three Options to Consider for the Minimum Number of Coins for Eleven Cents

Listing 8 is a dynamic programming algorithm to solve our change-making problem. `dpMakeChange` takes three parameters: a list of valid coin values, the amount of change we want to make, and a list of the minimum number of coins needed to make each value. When the function is done `minCoins` will contain the solution for all values from 0 to the value of change.

Listing 8

```
def dpMakeChange(coinValueList,change,minCoins):
    for cents in range(change+1):
        coinCount = cents
```

```

    for j in [c for c in coinValueList if c <= cents]:
        if minCoins[cents-j] + 1 < coinCount:
            coinCount = minCoins[cents-j]+1
    minCoins[cents] = coinCount
return minCoins[change]

```

Note that `dpMakeChange` is not a recursive function, even though we started with a recursive solution to this problem. It is important to realize that just because you can write a recursive solution to a problem does not mean it is the best or most efficient solution. The bulk of the work in this function is done by the loop that starts on line 4. In this loop we consider using all possible coins to make change for the amount specified by `cents`. Like we did for the 11 cent example above, we remember the minimum value and store it in our `minCoins` list.

Although our making change algorithm does a good job of figuring out the minimum number of coins, it does not help us make change since we do not keep track of the coins we use. We can easily extend `dpMakeChange` to keep track of the coins used by simply remembering the last coin we add for each entry in the `minCoins` table. If we know the last coin added, we can simply subtract the value of the coin to find a previous entry in the table that tells us the last coin we added to make that amount. We can keep tracing back through the table until we get to the beginning.

ActiveCode 2 shows the `dpMakeChange` algorithm modified to keep track of the coins used, along with a function `printCoins` that walks backward through the table to print out the value of each coin used. This shows the algorithm in action solving the problem for our friends in Lower Elbonia. The first two lines of `main` set the amount to be converted and create the list of coins used. The next two lines create the lists we need to store the results. `coinsUsed` is a list of the coins used to make change, and `coinCount` is the minimum number of coins used to make change for the amount corresponding to the position in the list.

Notice that the coins we print out come directly from the `coinsUsed` array. For the first call we start at array position 63 and print 21. Then we take  $63-21=42$  and look at the 42nd element of the list. Once again we find a 21 stored there. Finally, element 21 of the array also contains 21, giving us the three 21 cent pieces.

```

def dpMakeChange(coinValueList,change,minCoins,coinsUsed):
    for cents in range(change+1):
        coinCount = cents
        newCoin = 1
        for j in [c for c in coinValueList if c <= cents]:
            if minCoins[cents-j] + 1 < coinCount:
                coinCount = minCoins[cents-j]+1
                newCoin = j
        minCoins[cents] = coinCount
        coinsUsed[cents] = newCoin

```

```

    return minCoins[change]

def printCoins(coinsUsed,change):
    coin = change
    while coin > 0:
        thisCoin = coinsUsed[coin]
        print(thisCoin)
        coin = coin - thisCoin

def main():
    amnt = 63
    clist = [1,5,10,21,25]
    coinsUsed = [0]*(amnt+1)
    coinCount = [0]*(amnt+1)

    print("Making change for",amnt,"requires")
    print(dpMakeChange(clist,amnt,coinCount,coinsUsed),"coins")
    print("They are:")
    printCoins(coinsUsed,amnt)
    print("The used list is as follows:")
    print(coinsUsed)

main()

```

Complete Solution to the Change Problem (lst\_dpremember)

## 4.8. 小结

In this chapter we have looked at examples of several recursive algorithms. These algorithms were chosen to expose you to several different problems where recursion is an effective problem-solving technique. The key points to remember from this chapter are as follows:

- All recursive algorithms must have a base case.
- A recursive algorithm must change its state and make progress toward the base case.
- A recursive algorithm must call itself (recursively).
- Recursion can take the place of iteration in some cases.
- Recursive algorithms often map very naturally to a formal expression of the problem you are trying to solve.
- Recursion is not always the answer. Sometimes a recursive solution may be more computationally expensive than an alternative algorithm.

## 4.9. 关键词

base case	decrypt	dynamic programming
recursion	recursive call	stack frame

## 4.10. 问题讨论

1. Draw a call stack for the Tower of Hanoi problem. Assume that you start with a stack of three disks.
2. Using the recursive rules as described, draw a Sierpinski triangle using paper and pencil.
3. Using the dynamic programming algorithm for making change, find the smallest number of coins that you can use to make 33 cents in change. In addition to the usual coins assume that you have an 8 cent coin.

## 4.11. 词汇表

### **base case**

A branch of the conditional statement in a recursive function that does not give rise to further recursive calls.

### **data structure**

An organization of data for the purpose of making it easier to use.

### **exception**

An error that occurs at runtime.

### **handle an exception**

To prevent an exception from terminating a program by wrapping the block of code in a try / except construct.

### **immutable data type**

A data type which cannot be modified. Assignments to elements or slices of immutable types cause a runtime error.

### **infinite recursion**

A function that calls itself recursively without ever reaching the base case. Eventually, an infinite recursion causes a runtime error.

### **mutable data type**

A data type which can be modified. All mutable types are compound types. Lists and dictionaries (see next chapter) are mutable data types; strings and tuples are not.

### **raise**

To cause an exception by using the raise statement.

### **recursion**

The process of calling the function that is already executing.

### **recursive call**

The statement that calls an already executing function. Recursion can even be indirect —

function `f` can call `g` which calls `h`, and `h` could make a call back to `f`.

#### **recursive definition**

A definition which defines something in terms of itself. To be useful it must include base cases which are not recursive. In this way it differs from a circular definition. Recursive definitions often provide an elegant way to express complex data structures.

#### **tuple**

A data type that contains a sequence of elements of any type, like a list, but is immutable. Tuples can be used wherever an immutable type is required, such as a key in a dictionary (see next chapter).

#### **tuple assignment**

An assignment to all of the elements in a tuple using a single assignment statement. Tuple assignment occurs in parallel rather than in sequence, making it useful for swapping values.

## **4.12. 编程练习**

1. Write a recursive function to compute the factorial of a number.
2. Write a recursive function to reverse a list.
3. Modify the recursive tree program using one or all of the following ideas:
  - Modify the thickness of the branches so that as the `branchLen` gets smaller, the line gets thinner.
  - Modify the color of the branches so that as the `branchLen` gets very short it is colored like a leaf.
  - Modify the angle used in turning the turtle so that at each branch point the angle is selected at random in some range. For example choose the angle between 15 and 45 degrees. Play around to see what looks good.
  - Modify the `branchLen` recursively so that instead of always subtracting the same amount you subtract a random amount in some range.

If you implement all of the above ideas you will have a very realistic looking tree.

4. Find or invent an algorithm for drawing a fractal mountain. Hint: One approach to this uses triangles again.
5. Write a recursive function to compute the Fibonacci sequence. How does the performance of the recursive function compare to that of an iterative version?
6. Implement a solution to the Tower of Hanoi using three stacks to keep track of the disks.
7. Using the turtle graphics module, write a recursive program to display a Hilbert curve.
8. Using the turtle graphics module, write a recursive program to display a Koch snowflake.

9. Write a program to solve the following problem: You have two jugs: a 4-gallon jug and a 3-gallon jug. Neither of the jugs have markings on them. There is a pump that can be used to fill the jugs with water. How can you get exactly two gallons of water in the 4-gallon jug?
10. Generalize the problem above so that the parameters to your solution include the sizes of each jug and the final amount of water to be left in the larger jug.
11. Write a program that solves the following problem: Three missionaries and three cannibals come to a river and find a boat that holds two people. Everyone must get across the river to continue on the journey. However, if the cannibals ever outnumber the missionaries on either bank, the missionaries will be eaten. Find a series of crossings that will get everyone safely to the other side of the river.
12. Modify the Tower of Hanoi program using turtle graphics to animate the movement of the disks. Hint: You can make multiple turtles and have them shaped like rectangles.
13. Pascal' s triangle is a number triangle with numbers arranged in staggered rows such that

$$a_{nr} = \frac{n!}{r!(n-r)!}$$

This equation is the equation for a binomial coefficient. You can build Pascal' s triangle by adding the two numbers that are diagonally above a number in the triangle. An example of Pascal' s triangle is shown below.

```

      1
     1 1
    1 2 1
   1 3 3 1
  1 4 6 4 1

```

Write a program that prints out Pascal' s triangle. Your program should accept a parameter that tells how many rows of the triangle to print.

14. Suppose you are a computer scientist/art thief who has broken into a major art gallery. All you have with you to haul out your stolen art is your knapsack which only holds W pounds of art, but for every piece of art you know its value and its weight. Write a dynamic programming function to help you maximize your profit. Here is a sample problem for you to use to get started: Suppose your knapsack can hold a total weight of 20. You have 5 items as follows:

item	weight	value
1	2	3

2	3	4
3	4	8
4	5	8
5	9	10

15. This problem is called the string edit distance problem, and is quite useful in many areas of research. Suppose that you want to transform the word “algorithm” into the word “alligator.” For each letter you can either copy the letter from one word to another at a cost of 5, you can delete a letter at cost of 20, or insert a letter at a cost of 20. The total cost to transform one word into another is used by spell check programs to provide suggestions for words that are close to one another. Use dynamic programming techniques to develop an algorithm that gives you the smallest edit distance between any two words.

## 5. 排序与搜索

### 5.1. 目标

- 了解和实现顺序搜索和二分法搜索；
- 了解和实现选择排序、冒泡排序、归并排序、快速排序、插入排序和希尔排序；
- 了解用散列 Hashing 实现搜索的技术；
- 了解抽象数据类型：映射 Map；
- 采用散列实现抽象数据类型 Map。

### 5.2. 搜索

#### 5.2.1. 顺序搜索

When data items are stored in a collection such as a list, we say that they have a linear or sequential relationship. Each data item is stored in a position relative to the others. In Python lists, these relative positions are the index values of the individual items. Since these index values are ordered, it is possible for us to visit them in sequence. This process gives rise to our first searching technique, the sequential search.

Figure 1 shows how this search works. Starting at the first item in the list, we simply move from item to item, following the underlying sequential ordering until we either find what we are looking for or run out of items. If we run out of items, we have discovered that the item we were searching for was not present.

../\_images/seqsearch.png



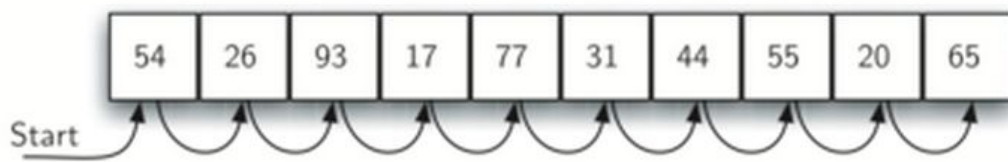


Figure 1: Sequential Search of a List of Integers

The Python implementation for this algorithm is shown in CodeLens 1. The function needs the list and the item we are looking for and returns a boolean value as to whether it is present. The boolean variable `found` is initialized to `False` and is assigned the value `True` if we discover the item in the list.

```
1 def sequentialSearch(alist, item):
2     pos = 0
3     found = False
4
5     while pos < len(alist) and not found:
6         if alist[pos] == item:
7             found = True
8         else:
9             pos = pos+1
10
11     return found
12
13 testlist = [1, 2, 32, 8, 17, 19, 42, 13, 0]
14 print(sequentialSearch(testlist, 3))
15 print(sequentialSearch(testlist, 13))
```

Sequential Search of an Unordered List (search1)

### 5.2.1.1. Analysis of Sequential Search

To analyze searching algorithms, we need to decide on a basic unit of computation. Recall that this is typically the common step that must be repeated in order to solve the problem. For searching, it makes sense to count the number of comparisons performed. Each comparison may or may not discover the item we are looking for. In addition, we make another assumption here. The list of items is not ordered in any way. The items have been placed randomly into the list. In other words, the probability that the item we are looking for is in any particular position is exactly the same for each position of the list.

If the item is not in the list, the only way to know it is to compare it against every item present. If there are  $n$  items, then the sequential search requires  $n$  comparisons to discover that the item is not

there. In the case where the item is in the list, the analysis is not so straightforward. There are actually three different scenarios that can occur. In the best case we will find the item in the first place we look, at the beginning of the list. We will need only one comparison. In the worst case, we will not discover the item until the very last comparison, the  $n$ th comparison.

What about the average case? On average, we will find the item about halfway into the list; that is, we will compare against  $n/2$  items. Recall, however, that as  $n$  gets large, the coefficients, no matter what they are, become insignificant in our approximation, so the complexity of the sequential search, is  $O(n)$ . Table 1 summarizes these results.

Table 1: Comparisons Used in a Sequential Search of an Unordered List

Case	Best Case	Worst Case	Average Case
item is present	1	$n$	$n/2$
item is not present	$n$	$n$	$n$

We assumed earlier that the items in our collection had been randomly placed so that there is no relative order between the items. What would happen to the sequential search if the items were ordered in some way? Would we be able to gain any efficiency in our search technique?

Assume that the list of items was constructed so that the items were in ascending order, from low to high. If the item we are looking for is present in the list, the chance of it being in any one of the  $n$  positions is still the same as before. We will still have the same number of comparisons to find the item. However, if the item is not present there is a slight advantage. Figure 2 shows this process as the algorithm looks for the item 50. Notice that items are still compared in sequence until 54. At this point, however, we know something extra. Not only is 54 not the item we are looking for, but no other elements beyond 54 can work either since the list is sorted. In this case, the algorithm does not have to continue looking through all of the items to report that the item was not found. It can stop immediately. CodeLens 2 shows this variation of the sequential search function.

../\_images/seqsearch2.png

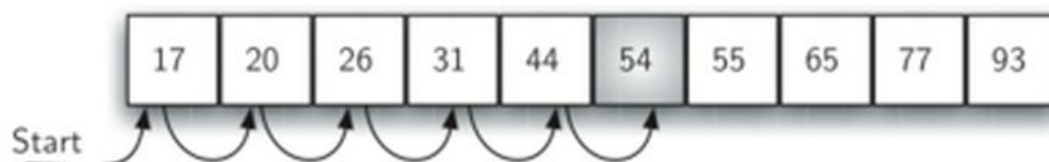


Figure 2: Sequential Search of an Ordered List of Integers

```

1  def orderedSequentialSearch(alist, item):
2      pos = 0
3      found = False
4      stop = False
5      while pos < len(alist) and not found and not stop:

```

```

6         if alist[pos] == item:
7             found = True
8         else:
9             if alist[pos] > item:
10                 stop = True
11             else:
12                 pos = pos+1
13
14     return found
15
16 testlist = [0, 1, 2, 8, 13, 17, 19, 32, 42,]
17 print(orderedSequentialSearch(testlist, 3))
18 print(orderedSequentialSearch(testlist, 13))

```

#### Sequential Search of an Ordered List (search2)

Table 2 summarizes these results. Note that in the best case we might discover that the item is not in the list by looking at only one item. On average, we will know after looking through only  $n/2$  items. However, this technique is still  $O(n)$ . In summary, a sequential search is improved by ordering the list only in the case where we do not find the item.

Table 2: Comparisons Used in Sequential Search of an Ordered List

Case	Best Case	Worst Case	Average Case
item is present	1	$n$	$n/2$
item not present	1	$n$	$n/2$

#### Self Check

Q-29: Suppose you are doing a sequential search of the list [15, 18, 2, 19, 18, 0, 8, 14, 19, 14]. How many comparisons would you need to do in order to find the key 18?

- a) 5
- b) 10
- c) 4
- d) 2

Q-30: Suppose you are doing a sequential search of the ordered list [3, 5, 6, 8, 11, 12, 14, 15, 17, 18]. How many comparisons would you need to do in order to find the key 13?

- a) 10
- b) 5
- c) 7
- d) 6

### 5.2.2. 二分法搜索

It is possible to take greater advantage of the ordered list if we are clever with our comparisons. In the sequential search, when we compare against the first item, there are at most  $n-1$  more items to look through if the first item is not what we are looking for. Instead of searching the list in sequence, a binary search will start by examining the middle item. If that item is the one we are searching for, we are done. If it is not the correct item, we can use the ordered nature of the list to eliminate half of the remaining items. If the item we are searching for is greater than the middle item, we know that the entire lower half of the list as well as the middle item can be eliminated from further consideration. The item, if it is in the list, must be in the upper half.

We can then repeat the process with the upper half. Start at the middle item and compare it against what we are looking for. Again, we either find it or split the list in half, therefore eliminating another large part of our possible search space. Figure 3 shows how this algorithm can quickly find the value 54. The complete function is shown in CodeLens 3.

../\_images/binsearch.png

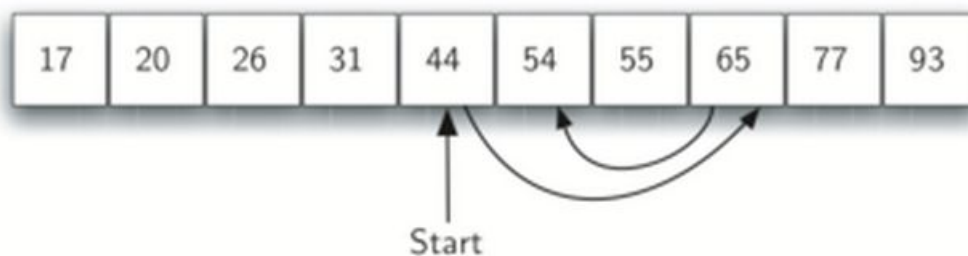


Figure 3: Binary Search of an Ordered List of Integers

```
1 def binarySearch(alist, item):
2     first = 0
3     last = len(alist)-1
4     found = False
5
6     while first<=last and not found:
7         midpoint = (first + last)//2
8         if alist[midpoint] == item:
9             found = True
10        else:
11            if item < alist[midpoint]:
12                last = midpoint-1
13            else:
14                first = midpoint+1
15
16    return found
```

```

17
18 testlist = [0, 1, 2, 8, 13, 17, 19, 32, 42,]
19 print(binarySearch(testlist, 3))
20 print(binarySearch(testlist, 13))

```

Binary Search of an Ordered List (search3)

Before we move on to the analysis, we should note that this algorithm is a great example of a divide and conquer strategy. Divide and conquer means that we divide the problem into smaller pieces, solve the smaller pieces in some way, and then reassemble the whole problem to get the result. When we perform a binary search of a list, we first check the middle item. If the item we are searching for is less than the middle item, we can simply perform a binary search of the left half of the original list. Likewise, if the item is greater, we can perform a binary search of the right half. Either way, this is a recursive call to the binary search function passing a smaller list. CodeLens 4 shows this recursive version.

```

1  def binarySearch(alist, item):
2      if len(alist) == 0:
3          return False
4      else:
5          midpoint = len(alist)//2
6          if alist[midpoint]==item:
7              return True
8          else:
9              if item<alist[midpoint]:
10                 return binarySearch(alist[:midpoint],item)
11             else:
12                 return binarySearch(alist[midpoint+1:],item)
13
14 testlist = [0, 1, 2, 8, 13, 17, 19, 32, 42,]
15 print(binarySearch(testlist, 3))
16 print(binarySearch(testlist, 13))

```

A Binary Search--Recursive Version (search4)

### 5.2.2.1. Analysis of Binary Search

To analyze the binary search algorithm, we need to recall that each comparison eliminates about half of the remaining items from consideration. What is the maximum number of comparisons this algorithm will require to check the entire list? If we start with  $n$  items, about  $n/2$  items will be left after the first comparison. After the second comparison, there will be about  $n/4$ . Then  $n/8$ ,  $n/16$ , and so on. How many times can we split the list? Table 3 helps us to see the answer.

Table 3: Tabular Analysis for a Binary Search

Comparisons	Approximate Number of Items Left
1	$n/2$
2	$n/4$
3	$n/8$
...	
$i$	$n/2^i$

When we split the list enough times, we end up with a list that has just one item. Either that is the item we are looking for or it is not. Either way, we are done. The number of comparisons necessary to get to this point is  $i$  where  $n/2^i=1$ . Solving for  $i$  gives us  $i=\log n$ . The maximum number of comparisons is logarithmic with respect to the number of items in the list. Therefore, the binary search is  $O(\log n)$ .

One additional analysis issue needs to be addressed. In the recursive solution shown above, the recursive call,

```
binarySearch(alist[:midpoint], item)
```

uses the slice operator to create the left half of the list that is then passed to the next invocation (similarly for the right half as well). The analysis that we did above assumed that the slice operator takes constant time. However, we know that the slice operator in Python is actually  $O(k)$ . This means that the binary search using slice will not perform in strict logarithmic time. Luckily this can be remedied by passing the list along with the starting and ending indices. The indices can be calculated as we did in Listing 3. We leave this implementation as an exercise.

Even though a binary search is generally better than a sequential search, it is important to note that for small values of  $n$ , the additional cost of sorting is probably not worth it. In fact, we should always consider whether it is cost effective to take on the extra work of sorting to gain searching benefits. If we can sort once and then search many times, the cost of the sort is not so significant. However, for large lists, sorting even once can be so expensive that simply performing a sequential search from the start may be the best choice.

### Self Check

Q-19: Suppose you have the following sorted list [3, 5, 6, 8, 11, 12, 14, 15, 17, 18] and are using the recursive binary search algorithm. Which group of numbers correctly shows the sequence of comparisons used to find the key 8.

- a) 11, 5, 6, 8
- b) 12, 6, 11, 8
- c) 3, 5, 6, 8
- d) 18, 12, 6, 8

Q-20: Suppose you have the following sorted list [3, 5, 6, 8, 11, 12, 14, 15, 17, 18] and are using the recursive binary search algorithm. Which group of numbers correctly shows the sequence of

comparisons used to search for the key 16?

- a) 11, 14, 17
- b) 18, 17, 15
- c) 14, 17, 15
- d) 12, 17, 15

### 5.2.3. 散列

In previous sections we were able to make improvements in our search algorithms by taking advantage of information about where items are stored in the collection with respect to one another. For example, by knowing that a list was ordered, we could search in logarithmic time using a binary search. In this section we will attempt to go one step further by building a data structure that can be searched in  $O(1)$  time. This concept is referred to as hashing.

In order to do this, we will need to know even more about where the items might be when we go to look for them in the collection. If every item is where it should be, then the search can use a single comparison to discover the presence of an item. We will see, however, that this is typically not the case.

A hash table is a collection of items which are stored in such a way as to make it easy to find them later. Each position of the hash table, often called a slot, can hold an item and is named by an integer value starting at 0. For example, we will have a slot named 0, a slot named 1, a slot named 2, and so on. Initially, the hash table contains no items so every slot is empty. We can implement a hash table by using a list with each element initialized to the special Python value `None`. Figure 4 shows a hash table of size  $m=11$ . In other words, there are  $m$  slots in the table, named 0 through 10.

../\_images/hashtable.png

0	1	2	3	4	5	6	7	8	9	10
None	None	None	None	None	None	None	None	None	None	None

Figure 4: Hash Table with 11 Empty Slots

The mapping between an item and the slot where that item belongs in the hash table is called the hash function. The hash function will take any item in the collection and return an integer in the range of slot names, between 0 and  $m-1$ . Assume that we have the set of integer items 54, 26, 93, 17, 77, and 31. Our first hash function, sometimes referred to as the “remainder method,” simply takes an item and divides it by the table size, returning the remainder as its hash value ( $h(\text{item}) = \text{item} \% 11$ ). Table 4 gives all of the hash values for our example items. Note that this remainder method (modulo arithmetic) will typically be present in some form in all hash functions, since the result must be in the range of slot names.

Table 4: Simple Hash Function Using Remainders

Item	Hash Value
54	10
26	4
93	5
17	6
77	0
31	9

Once the hash values have been computed, we can insert each item into the hash table at the designated position as shown in Figure 5. Note that 6 of the 11 slots are now occupied. This is referred to as the load factor, and is commonly denoted by  $\lambda = \text{number of items} / \text{table size}$ . For this example,  $\lambda = 6/11$ .

../\_images/hashtable2.png

0	1	2	3	4	5	6	7	8	9	10
77	None	None	None	26	93	17	None	None	31	54

Figure 5: Hash Table with Six Items

Now when we want to search for an item, we simply use the hash function to compute the slot name for the item and then check the hash table to see if it is present. This searching operation is  $O(1)$ , since a constant amount of time is required to compute the hash value and then index the hash table at that location. If everything is where it should be, we have found a constant time search algorithm.

You can probably already see that this technique is going to work only if each item maps to a unique location in the hash table. For example, if the item 44 had been the next item in our collection, it would have a hash value of 0 ( $44 \% 11 = 0$ ). Since 77 also had a hash value of 0, we would have a problem. According to the hash function, two or more items would need to be in the same slot. This is referred to as a collision (it may also be called a “clash”). Clearly, collisions create a problem for the hashing technique. We will discuss them in detail later.

### 5.2.3.1. Hash Functions

Given a collection of items, a hash function that maps each item into a unique slot is referred to as a perfect hash function. If we know the items and the collection will never change, then it is possible to construct a perfect hash function (refer to the exercises for more about perfect hash functions). Unfortunately, given an arbitrary collection of items, there is no systematic way to construct a perfect hash function. Luckily, we do not need the hash function to be perfect to still gain performance efficiency.



One way to always have a perfect hash function is to increase the size of the hash table so that each possible value in the item range can be accommodated. This guarantees that each item will have a unique slot. Although this is practical for small numbers of items, it is not feasible when the number of possible items is large. For example, if the items were nine-digit Social Security numbers, this method would require almost one billion slots. If we only want to store data for a class of 25 students, we will be wasting an enormous amount of memory.

Our goal is to create a hash function that minimizes the number of collisions, is easy to compute, and evenly distributes the items in the hash table. There are a number of common ways to extend the simple remainder method. We will consider a few of them here.

The folding method for constructing hash functions begins by dividing the item into equal-size pieces (the last piece may not be of equal size). These pieces are then added together to give the resulting hash value. For example, if our item was the phone number 436-555-4601, we would take the digits and divide them into groups of 2 (43,65,55,46,01). After the addition,  $43+65+55+46+01$ , we get 210. If we assume our hash table has 11 slots, then we need to perform the extra step of dividing by 11 and keeping the remainder. In this case  $210 \% 11$  is 1, so the phone number 436-555-4601 hashes to slot 1. Some folding methods go one step further and reverse every other piece before the addition. For the above example, we get  $43+56+55+64+01=219$  which gives  $219 \% 11=10$ .

Another numerical technique for constructing a hash function is called the mid-square method. We first square the item, and then extract some portion of the resulting digits. For example, if the item were 44, we would first compute  $44^2=1,936$ . By extracting the middle two digits, 93, and performing the remainder step, we get 5 ( $93 \% 11$ ). Table 5 shows items under both the remainder method and the mid-square method. You should verify that you understand how these values were computed.

Table 5: Comparison of Remainder and Mid-Square Methods

Item	Remainder	Mid-Square
54	10	3
26	4	7
93	5	9
17	6	8
77	0	4
31	9	6

We can also create hash functions for character-based items such as strings. The word “cat” can be thought of as a sequence of ordinal values.

```
>>> ord('c')
99
>>> ord('a')
97
>>> ord('t')
```

We can then take these three ordinal values, add them up, and use the remainder method to get a hash value (see Figure 6). Listing 1 shows a function called `hash` that takes a string and a table size and returns the hash value in the range from 0 to `tablesize-1`.

../\_images/stringhash.png

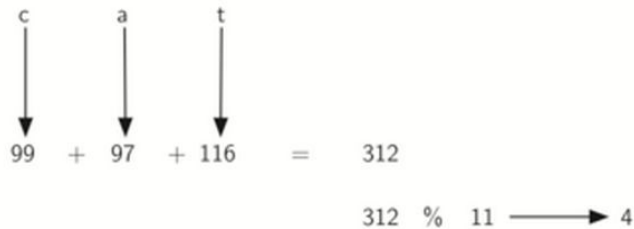


Figure 6: Hashing a String Using Ordinal Values

Listing 1

```
def hash(astring, tablesize):
    sum = 0
    for pos in range(len(astring)):
        sum = sum + ord(astring[pos])

    return sum%tablesize
```

It is interesting to note that when using this hash function, anagrams will always be given the same hash value. To remedy this, we could use the position of the character as a weight. Figure 7 shows one possible way to use the positional value as a weighting factor. The modification to the hash function is left as an exercise.

../\_images/stringhash2.png

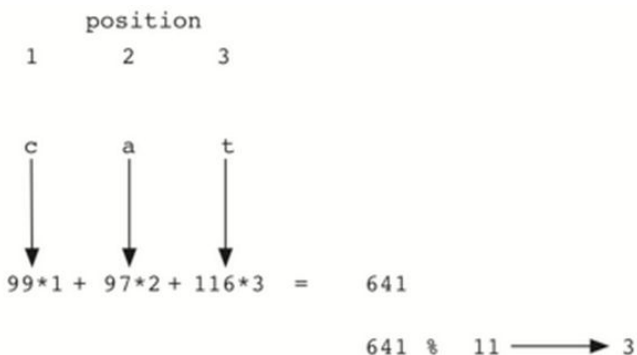


Figure 7: Hashing a String Using Ordinal Values with Weighting

You may be able to think of a number of additional ways to compute hash values for items in a collection. The important thing to remember is that the hash function has to be efficient so that it does not become the dominant part of the storage and search process. If the hash function is too

complex, then it becomes more work to compute the slot name than it would be to simply do a basic sequential or binary search as described earlier. This would quickly defeat the purpose of hashing.

### 5.2.3.2. Collision Resolution

We now return to the problem of collisions. When two items hash to the same slot, we must have a systematic method for placing the second item in the hash table. This process is called collision resolution. As we stated earlier, if the hash function is perfect, collisions will never occur. However, since this is often not possible, collision resolution becomes a very important part of hashing.

One method for resolving collisions looks into the hash table and tries to find another open slot to hold the item that caused the collision. A simple way to do this is to start at the original hash value position and then move in a sequential manner through the slots until we encounter the first slot that is empty. Note that we may need to go back to the first slot (circularly) to cover the entire hash table. This collision resolution process is referred to as open addressing in that it tries to find the next open slot or address in the hash table. By systematically visiting each slot one at a time, we are performing an open addressing technique called linear probing.

Figure 8 shows an extended set of integer items under the simple remainder method hash function (54,26,93,17,77,31,44,55,20). Table 4 above shows the hash values for the original items. Figure 5 shows the original contents. When we attempt to place 44 into slot 0, a collision occurs. Under linear probing, we look sequentially, slot by slot, until we find an open position. In this case, we find slot 1.

Again, 55 should go in slot 0 but must be placed in slot 2 since it is the next open position. The final value of 20 hashes to slot 9. Since slot 9 is full, we begin to do linear probing. We visit slots 10, 0, 1, and 2, and finally find an empty slot at position 3.

../\_images/linearprobing1.png

0	1	2	3	4	5	6	7	8	9	10
77	44	55	20	26	93	17	None	None	31	54

Figure 8: Collision Resolution with Linear Probing

Once we have built a hash table using open addressing and linear probing, it is essential that we utilize the same methods to search for items. Assume we want to look up the item 93. When we compute the hash value, we get 5. Looking in slot 5 reveals 93, and we can return True. What if we are looking for 20? Now the hash value is 9, and slot 9 is currently holding 31. We cannot simply return False since we know that there could have been collisions. We are now forced to do a sequential search, starting at position 10, looking until either we find the item 20 or we find an

empty slot.

A disadvantage to linear probing is the tendency for clustering; items become clustered in the table. This means that if many collisions occur at the same hash value, a number of surrounding slots will be filled by the linear probing resolution. This will have an impact on other items that are being inserted, as we saw when we tried to add the item 20 above. A cluster of values hashing to 0 had to be skipped to finally find an open position. This cluster is shown in Figure 9.

../\_images/clustering.png

0	1	2	3	4	5	6	7	8	9	10
77	44	55	20	26	93	17	None	None	31	54

Figure 9: A Cluster of Items for Slot 0

One way to deal with clustering is to extend the linear probing technique so that instead of looking sequentially for the next open slot, we skip slots, thereby more evenly distributing the items that have caused collisions. This will potentially reduce the clustering that occurs. Figure 10 shows the items when collision resolution is done with a “plus 3” probe. This means that once a collision occurs, we will look at every third slot until we find one that is empty.

../\_images/linearprobing2.png

0	1	2	3	4	5	6	7	8	9	10
77	55	None	44	26	93	17	20	None	31	54

Figure 10: Collision Resolution Using “Plus 3”

The general name for this process of looking for another slot after a collision is rehashing. With simple linear probing, the rehash function is  $\text{newhashvalue} = \text{rehash}(\text{oldhashvalue})$  where  $\text{rehash}(\text{pos}) = (\text{pos} + 1) \% \text{sizeof table}$ . The “plus 3” rehash can be defined as  $\text{rehash}(\text{pos}) = (\text{pos} + 3) \% \text{sizeof table}$ . In general,  $\text{rehash}(\text{pos}) = (\text{pos} + \text{skip}) \% \text{sizeof table}$ . It is important to note that the size of the “skip” must be such that all the slots in the table will eventually be visited. Otherwise, part of the table will be unused. To ensure this, it is often suggested that the table size be a prime number. This is the reason we have been using 11 in our examples.

A variation of the linear probing idea is called quadratic probing. Instead of using a constant “skip” value, we use a rehash function that increments the hash value by 1, 3, 5, 7, 9, and so on. This means that if the first hash value is  $h$ , the successive values are  $h+1$ ,  $h+4$ ,  $h+9$ ,  $h+16$ , and so on. In other words, quadratic probing uses a skip consisting of successive perfect squares. Figure 11 shows our example values after they are placed using this technique.

../\_images/quadratic.png

0	1	2	3	4	5	6	7	8	9	10
77	44	20	55	26	93	17	None	None	31	54

Figure 11: Collision Resolution with Quadratic Probing

An alternative method for handling the collision problem is to allow each slot to hold a reference to a collection (or chain) of items. Chaining allows many items to exist at the same location in the hash table. When collisions happen, the item is still placed in the proper slot of the hash table. As more and more items hash to the same location, the difficulty of searching for the item in the collection increases. Figure 12 shows the items as they are added to a hash table that uses chaining to resolve collisions.

../\_images/chaining.png

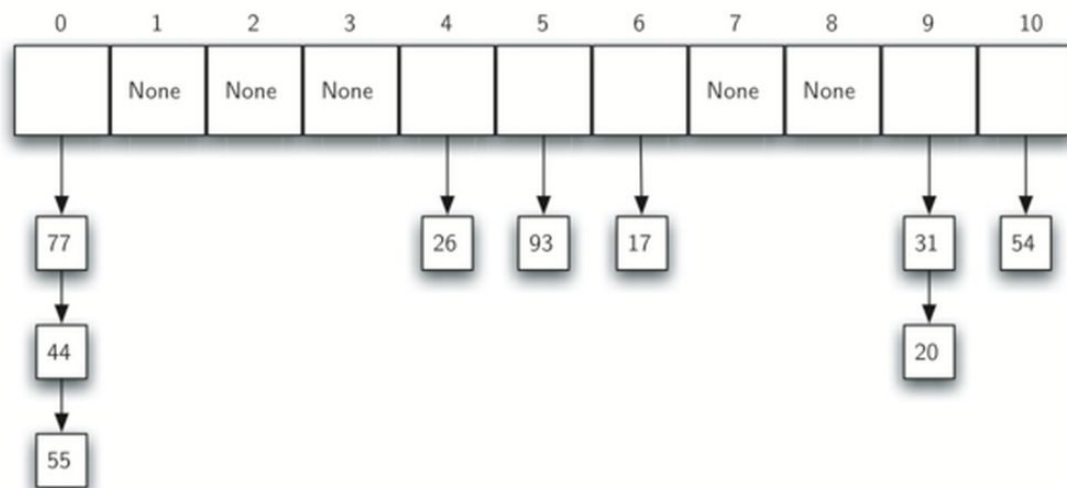


Figure 12: Collision Resolution with Chaining

When we want to search for an item, we use the hash function to generate the slot where it should reside. Since each slot holds a collection, we use a searching technique to decide whether the item is present. The advantage is that on the average there are likely to be many fewer items in each slot, so the search is perhaps more efficient. We will look at the analysis for hashing at the end of this section.

### Self Check

Q-17: In a hash table of size 13 which index positions would the following two keys map to? 27, 130

- a) 1, 10
- b) 13, 0
- c) 1, 0
- d) 2, 3

Q-18: Suppose you are given the following set of keys to insert into a hash table that holds exactly

11 values: 113 , 117 , 97 , 100 , 114 , 108 , 116 , 105 , 99 Which of the following best demonstrates the contents of the has table after all the keys have been inserted using linear probing?

- a) 100, \_\_, \_\_, 113, 114, 105, 116, 117, 97, 108, 99
- b) 99, 100, \_\_, 113, 114, \_\_, 116, 117, 105, 97, 108
- c) 100, 113, 117, 97, 14, 108, 116, 105, 99, \_\_, \_\_
- d) 117, 114, 108, 116, 105, 99, \_\_, \_\_, 97, 100, 113

### 5.2.3.3. Implementing the Map Abstract Data Type

One of the most useful Python collections is the dictionary. Recall that a dictionary is an associative data type where you can store key – data pairs. The key is used to look up the associated data value. We often refer to this idea as a map.

The map abstract data type is defined as follows. The structure is an unordered collection of associations between a key and a data value. The keys in a map are all unique so that there is a one-to-one relationship between a key and a value. The operations are given below.

- Map() Create a new, empty map. It returns an empty map collection.
- put(key,val) Add a new key-value pair to the map. If the key is already in the map then replace the old value with the new value.
- get(key) Given a key, return the value stored in the map or None otherwise.
- del Delete the key-value pair from the map using a statement of the form del map[key].
- len() Return the number of key-value pairs stored in the map.
- in Return True for a statement of the form key in map, if the given key is in the map, False otherwise.

One of the great benefits of a dictionary is the fact that given a key, we can look up the associated data value very quickly. In order to provide this fast look up capability, we need an implementation that supports an efficient search. We could use a list with sequential or binary search but it would be even better to use a hash table as described above since looking up an item in a hash table can approach  $O(1)$  performance.

In Listing 2 we use two lists to create a HashTable class that implements the Map abstract data type. One list, called slots, will hold the key items and a parallel list, called data, will hold the data values. When we look up a key, the corresponding position in the data list will hold the associated data value. We will treat the key list as a hash table using the ideas presented earlier. Note that the initial size for the hash table has been chosen to be 11. Although this is arbitrary, it is important that the size be a prime number so that the collision resolution algorithm can be as efficient as possible.

Listing 2

```

class HashTable:
    def __init__(self):
        self.size = 11
        self.slots = [None] * self.size
        self.data = [None] * self.size

```

hashfunction implements the simple remainder method. The collision resolution technique is linear probing with a “plus 1” rehash function. The put function (see Listing 3) assumes that there will eventually be an empty slot unless the key is already present in the self.slots. It computes the original hash value and if that slot is not empty, iterates the rehash function until an empty slot occurs. If a nonempty slot already contains the key, the old data value is replaced with the new data value.

Listing 3

```

def put(self, key, data):
    hashvalue = self.hashfunction(key, len(self.slots))

    if self.slots[hashvalue] == None:
        self.slots[hashvalue] = key
        self.data[hashvalue] = data
    else:
        if self.slots[hashvalue] == key:
            self.data[hashvalue] = data #replace
        else:
            nextslot = self.rehash(hashvalue, len(self.slots))
            while self.slots[nextslot] != None and \
                  self.slots[nextslot] != key:
                nextslot = self.rehash(nextslot, len(self.slots))

            if self.slots[nextslot] == None:
                self.slots[nextslot] = key
                self.data[nextslot] = data
            else:
                self.data[nextslot] = data #replace

def hashfunction(self, key, size):
    return key % size

def rehash(self, oldhash, size):
    return (oldhash + 1) % size

```

Likewise, the get function (see Listing 4) begins by computing the initial hash value. If the value is not in the initial slot, rehash is used to locate the next possible position. Notice that line 15 guarantees that the search will terminate by checking to make sure that we have not returned to the

initial slot. If that happens, we have exhausted all possible slots and the item must not be present.

The final methods of the HashTable class provide additional dictionary functionality. We overload the `__getitem__` and `__setitem__` methods to allow access using `[]`. This means that once a HashTable has been created, the familiar index operator will be available. We leave the remaining methods as exercises.

#### Listing 4

```
def get(self, key):
    startslot = self.hashfunction(key, len(self.slots))

    data = None
    stop = False
    found = False
    position = startslot
    while self.slots[position] != None and \
           not found and not stop:
        if self.slots[position] == key:
            found = True
            data = self.data[position]
        else:
            position = self.rehash(position, len(self.slots))
            if position == startslot:
                stop = True
    return data

def __getitem__(self, key):
    return self.get(key)

def __setitem__(self, key, data):
    self.put(key, data)
```

The following session shows the HashTable class in action. First we will create a hash table and store some items with integer keys and string data values.

```
>>> H=HashTable()
>>> H[54]="cat"
>>> H[26]="dog"
>>> H[93]="lion"
>>> H[17]="tiger"
>>> H[77]="bird"
>>> H[31]="cow"
>>> H[44]="goat"
```



```

>>> H[55]="pig"
>>> H[20]="chicken"
>>> H.slots
[77, 44, 55, 20, 26, 93, 17, None, None, 31, 54]
>>> H.data
['bird', 'goat', 'pig', 'chicken', 'dog', 'lion',
 'tiger', None, None, 'cow', 'cat']

```

Next we will access and modify some items in the hash table. Note that the value for the key 20 is being replaced.

```

>>> H[20]
'chicken'
>>> H[17]
'tiger'
>>> H[20]='duck'
>>> H[20]
'duck'
>>> H.data
['bird', 'goat', 'pig', 'duck', 'dog', 'lion',
 'tiger', None, None, 'cow', 'cat']
>> print(H[99])
None

```

The complete hash table example can be found in ActiveCode 1.

Complete Hash Table Example (hashtablecomplete)

### 5.2.3.4. Analysis of Hashing

We stated earlier that in the best case hashing would provide a  $O(1)$ , constant time search technique. However, due to collisions, the number of comparisons is typically not so simple. Even though a complete analysis of hashing is beyond the scope of this text, we can state some well-known results that approximate the number of comparisons necessary to search for an item.

The most important piece of information we need to analyze the use of a hash table is the load factor,  $\lambda$ . Conceptually, if  $\lambda$  is small, then there is a lower chance of collisions, meaning that items are more likely to be in the slots where they belong. If  $\lambda$  is large, meaning that the table is filling up, then there are more and more collisions. This means that collision resolution is more difficult, requiring more comparisons to find an empty slot. With chaining, increased collisions means an increased number of items on each chain.

As before, we will have a result for both a successful and an unsuccessful search. For a successful search using open addressing with linear probing, the average number of comparisons is

approximately  $12(1+1/1-\lambda)$  and an unsuccessful search gives  $12(1+(1/1-\lambda)^2)$  If we are using chaining, the average number of comparisons is  $1+\lambda/2$  for the successful case, and simply  $\lambda$  comparisons if the search is unsuccessful.

## 5.3. 排序

### 5.3.1. 冒泡排序

The bubble sort makes multiple passes through a list. It compares adjacent items and exchanges those that are out of order. Each pass through the list places the next largest value in its proper place. In essence, each item “bubbles” up to the location where it belongs.

Figure 1 shows the first pass of a bubble sort. The shaded items are being compared to see if they are out of order. If there are  $n$  items in the list, then there are  $n-1$  pairs of items that need to be compared on the first pass. It is important to note that once the largest value in the list is part of a pair, it will continually be moved along until the pass is complete.

../\_images/bubblepass.png

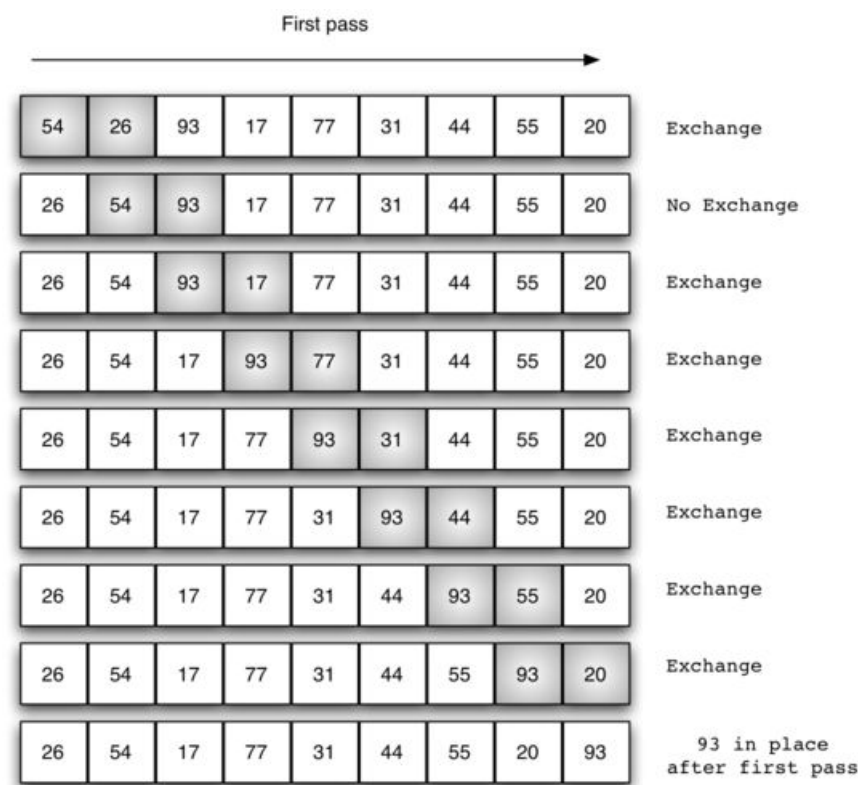


Figure 1: bubbleSort: The First Pass

At the start of the second pass, the largest value is now in place. There are  $n-1$  items left to sort, meaning that there will be  $n-2$  pairs. Since each pass places the next largest value in place, the total number of passes necessary will be  $n-1$ . After completing the  $n-1$  passes, the smallest item

must be in the correct position with no further processing required. ActiveCode 1 shows the complete bubbleSort function. It takes the list as a parameter, and modifies it by exchanging items as necessary.

The exchange operation, sometimes called a “swap,” is slightly different in Python than in most other programming languages. Typically, swapping two elements in a list requires a temporary storage location (an additional memory location). A code fragment such as

```
temp = alist[i]
alist[i] = alist[j]
alist[j] = temp
```

will exchange the *i*th and *j*th items in the list. Without the temporary storage, one of the values would be overwritten.

In Python, it is possible to perform simultaneous assignment. The statement `a,b=b,a` will result in two assignment statements being done at the same time (see Figure 2). Using simultaneous assignment, the exchange operation can be done in one statement.

Lines 5-7 in ActiveCode 1 perform the exchange of the *i* and (*i*+1)th items using the three – step procedure described earlier. Note that we could also have used the simultaneous assignment to swap the items.

../\_images/swap.png

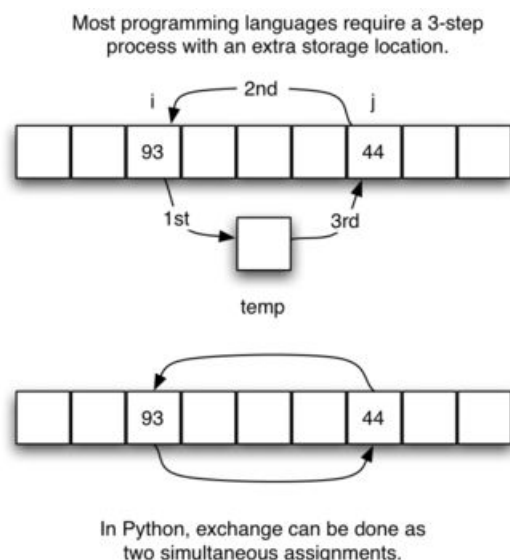


Figure 2: Exchanging Two Values in Python

The following activecode example shows the complete bubbleSort function working on the list shown above.

Run Save Load Show in Codelens

```

def bubbleSort(alist):
    for passnum in range(len(alist)-1,0,-1):
        for i in range(passnum):
            if alist[i]>alist[i+1]:
                temp = alist[i]
                alist[i] = alist[i+1]
                alist[i+1] = temp

alist = [54,26,93,17,77,31,44,55,20]
bubbleSort(alist)
print(alist)

```

The Bubble Sort (lst\_bubble)

The following animation shows bubbleSort in action.

Initialize Run Stop

Beginning Step Forward Step Backward End

To analyze the bubble sort, we should note that regardless of how the items are arranged in the initial list,  $n - 1$  passes will be made to sort a list of size  $n$ . Table 1 shows the number of comparisons for each pass. The total number of comparisons is the sum of the first  $n - 1$  integers. Recall that the sum of the first  $n$  integers is  $\frac{1}{2}n^2 + \frac{1}{2}n$ . The sum of the first  $n - 1$  integers is  $\frac{1}{2}n^2 + \frac{1}{2}n - n$ , which is  $\frac{1}{2}n^2 - \frac{1}{2}n$ . This is still  $O(n^2)$  comparisons. In the best case, if the list is already ordered, no exchanges will be made. However, in the worst case, every comparison will cause an exchange. On average, we exchange half of the time.

Table 1: Comparisons for Each Pass of Bubble Sort

Pass	Comparisons
1	$n-1$
2	$n-2$
3	$n-3$
...	...
$n-1$	1

A bubble sort is often considered the most inefficient sorting method since it must exchange items before the final location is known. These “wasted” exchange operations are very costly. However, because the bubble sort makes passes through the entire unsorted portion of the list, it has the capability to do something most sorting algorithms cannot. In particular, if during a pass there are no exchanges, then we know that the list must be sorted. A bubble sort can be modified to stop early if it finds that the list has become sorted. This means that for lists that require just a few passes, a bubble sort may have an advantage in that it will recognize the sorted list and stop. ActiveCode 2 shows this modification, which is often referred to as the short bubble.

Run Save Load Show in Codelens

```
def shortBubbleSort(alist):
    exchanges = True
    passnum = len(alist)-1
    while passnum > 0 and exchanges:
        exchanges = False
        for i in range(passnum):
            if alist[i]>alist[i+1]:
                exchanges = True
                temp = alist[i]
                alist[i] = alist[i+1]
                alist[i+1] = temp
        passnum = passnum-1

alist=[20,30,40,90,50,60,70,80,100,110]
shortBubbleSort(alist)
print(alist)
```

The Short Bubble Sort (lst\_shortbubble)

#### Self Check

Q-21: Suppose you have the following list of numbers to sort: [19, 1, 9, 7, 3, 10, 13, 15, 8, 12] which list represents the partially sorted list after three complete passes of bubble sort?

- a) [1, 9, 19, 7, 3, 10, 13, 15, 8, 12]
- b) [1, 3, 7, 9, 10, 8, 12, 13, 15, 19]
- c) [1, 7, 3, 9, 10, 13, 8, 12, 15, 19]
- d) [1, 9, 19, 7, 3, 10, 13, 15, 8, 12]

### 5.3.2. 选择排序

The selection sort improves on the bubble sort by making only one exchange for every pass through the list. In order to do this, a selection sort looks for the largest value as it makes a pass and, after completing the pass, places it in the proper location. As with a bubble sort, after the first pass, the largest item is in the correct place. After the second pass, the next largest is in place. This process continues and requires  $n-1$  passes to sort  $n$  items, since the final item must be in place after the  $(n-1)$  st pass.

Figure 3 shows the entire sorting process. On each pass, the largest remaining item is selected and then placed in its proper location. The first pass places 93, the second pass places 77, the third places 55, and so on. The function is shown in ActiveCode 1.

../\_images/selectionsortnew.png

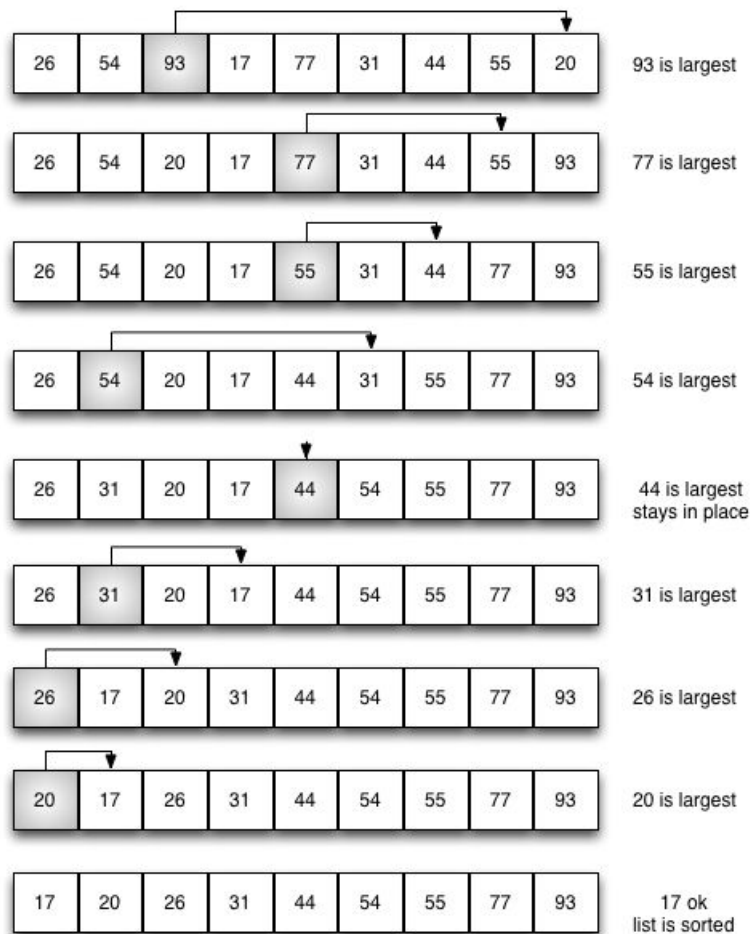


Figure 3: selectionSort

```
def selectionSort(alist):
    for fillslot in range(len(alist)-1,0,-1):
        positionOfMax=0
        for location in range(1,fillslot+1):
            if alist[location]>alist[positionOfMax]:
                positionOfMax = location

        temp = alist[fillslot]
        alist[fillslot] = alist[positionOfMax]
        alist[positionOfMax] = temp

alist = [54,26,93,17,77,31,44,55,20]
selectionSort(alist)
print(alist)
```

Selection Sort (lst\_selectionsortcode)

Initialize Run Stop

Beginning Step Forward Step Backward End

You may see that the selection sort makes the same number of comparisons as the bubble sort and is therefore also  $O(n^2)$ . However, due to the reduction in the number of exchanges, the selection sort typically executes faster in benchmark studies. In fact, for our list, the bubble sort makes 20 exchanges, while the selection sort makes only 8.

### Self Check

Q-28: Suppose you have the following list of numbers to sort: [11, 7, 12, 14, 19, 1, 6, 18, 8, 20] which list represents the partially sorted list after three complete passes of selection sort?

- a) [7, 11, 12, 1, 6, 14, 8, 18, 19, 20]
- b) [7, 11, 12, 14, 19, 1, 6, 18, 8, 20]
- c) [11, 7, 12, 14, 1, 6, 8, 18, 19, 20]
- d) [11, 7, 12, 14, 8, 1, 6, 18, 19, 20]

### 5.3.3. 插入排序

The insertion sort, although still  $O(n^2)$ , works in a slightly different way. It always maintains a sorted sublist in the lower positions of the list. Each new item is then “inserted” back into the previous sublist such that the sorted sublist is one item larger. Figure 4 shows the insertion sorting process. The shaded items represent the ordered sublists as the algorithm makes each pass.

../\_images/insertionsort.png

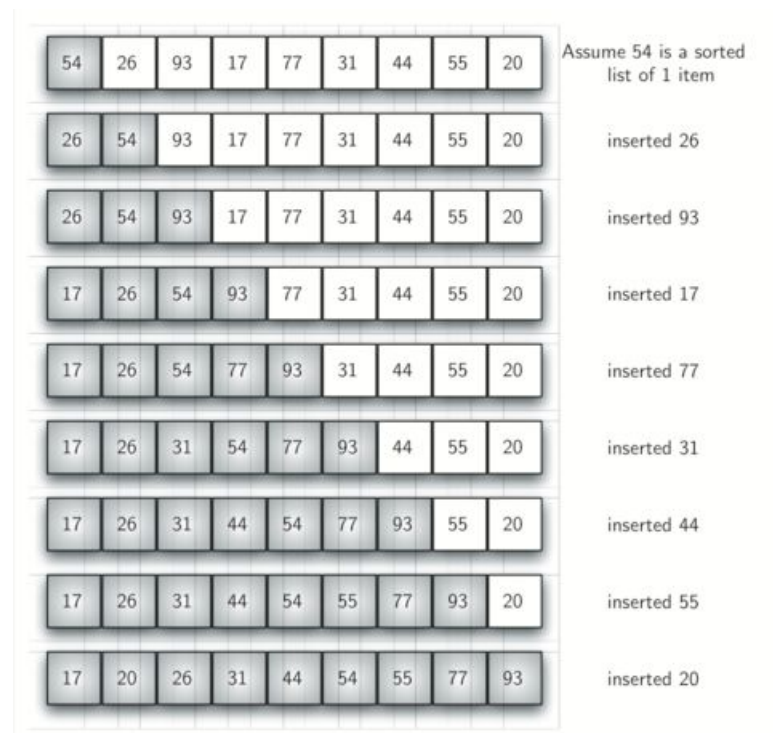


Figure 4: insertionSort

We begin by assuming that a list with one item (position 0) is already sorted. On each pass, one for each item 1 through  $n-1$ , the current item is checked against those in the already sorted sublist. As we look back into the already sorted sublist, we shift those items that are greater to the right. When we reach a smaller item or the end of the sublist, the current item can be inserted.

Figure 5 shows the fifth pass in detail. At this point in the algorithm, a sorted sublist of five items consisting of 17, 26, 54, 77, and 93 exists. We want to insert 31 back into the already sorted items. The first comparison against 93 causes 93 to be shifted to the right. 77 and 54 are also shifted. When the item 26 is encountered, the shifting process stops and 31 is placed in the open position. Now we have a sorted sublist of six items.

../\_images/insertionpass.png

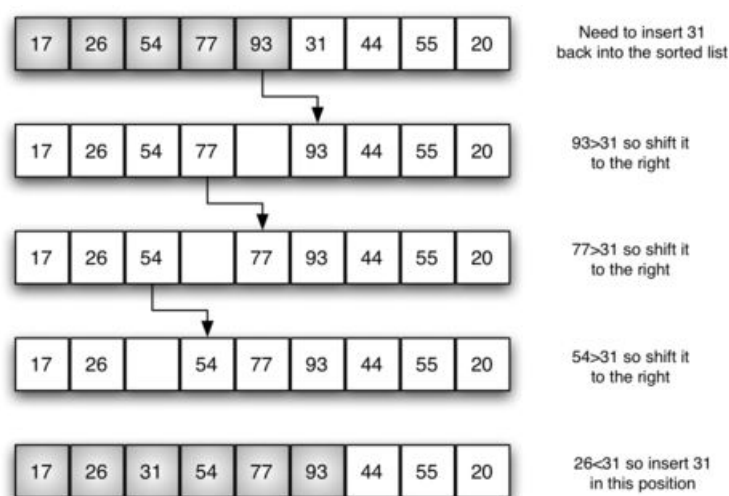


Figure 5: insertionSort: Fifth Pass of the Sort

The implementation of insertionSort (ActiveCode 1) shows that there are again  $n-1$  passes to sort  $n$  items. The iteration starts at position 1 and moves through position  $n-1$ , as these are the items that need to be inserted back into the sorted sublists. Line 8 performs the shift operation that moves a value up one position in the list, making room behind it for the insertion. Remember that this is not a complete exchange as was performed in the previous algorithms.

The maximum number of comparisons for an insertion sort is the sum of the first  $n-1$  integers. Again, this is  $O(n^2)$ . However, in the best case, only one comparison needs to be done on each pass. This would be the case for an already sorted list.

One note about shifting versus exchanging is also important. In general, a shift operation requires approximately a third of the processing work of an exchange since only one assignment is performed. In benchmark studies, insertion sort will show very good performance.

Run Save Load Show in Codelens



```

def insertionSort(alist):
    for index in range(1,len(alist)):

        currentvalue = alist[index]
        position = index

        while position>0 and alist[position-1]>currentvalue:
            alist[position]=alist[position-1]
            position = position-1

        alist[position]=currentvalue

alist = [54,26,93,17,77,31,44,55,20]
insertionSort(alist)
print(alist)

```

Insertion Sort (lst\_insertion)

Initialize Run Stop

Beginning Step Forward Step Backward End

#### Self Check

Q-22: Suppose you have the following list of numbers to sort: <br>

[15, 5, 4, 18, 12, 19, 14, 10, 8, 20] which list represents the partially sorted list after three complete passes of insertion sort?

- a) [4, 5, 12, 15, 14, 10, 8, 18, 19, 20]
- b) [15, 5, 4, 10, 12, 8, 14, 18, 19, 20]
- c) [4, 5, 15, 18, 12, 19, 14, 10, 8, 20]
- d) [15, 5, 4, 18, 12, 19, 14, 8, 10, 20]

### 5.3.4. 希尔排序

The shell sort, sometimes called the “diminishing increment sort,” improves on the insertion sort by breaking the original list into a number of smaller sublists, each of which is sorted using an insertion sort. The unique way that these sublists are chosen is the key to the shell sort. Instead of breaking the list into sublists of contiguous items, the shell sort uses an increment *i*, sometimes called the gap, to create a sublist by choosing all items that are *i* items apart.

This can be seen in Figure 6. This list has nine items. If we use an increment of three, there are three sublists, each of which can be sorted by an insertion sort. After completing these sorts, we get the list shown in Figure 7. Although this list is not completely sorted, something very interesting has happened. By sorting the sublists, we have moved the items closer to where they actually belong.

../\_images/shellsortA.png

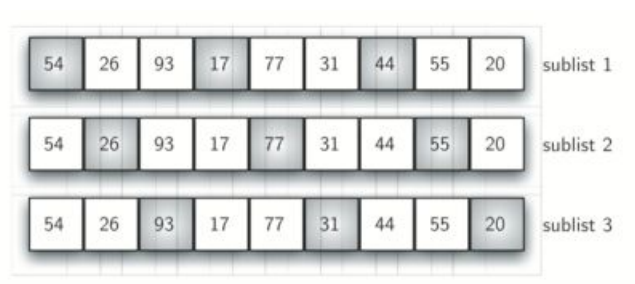


Figure 6: A Shell Sort with Increments of Three

../\_images/shellsortB.png

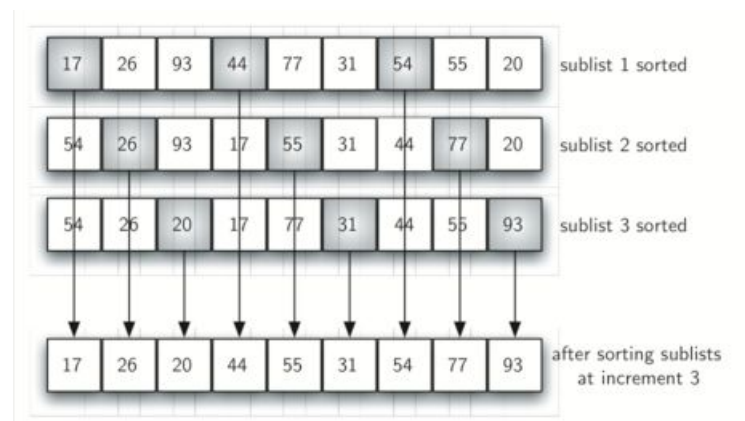


Figure 7: A Shell Sort after Sorting Each Sublist

Figure 8 shows a final insertion sort using an increment of one; in other words, a standard insertion sort. Note that by performing the earlier sublist sorts, we have now reduced the total number of shifting operations necessary to put the list in its final order. For this case, we need only four more shifts to complete the process.

../\_images/shellsortC.png

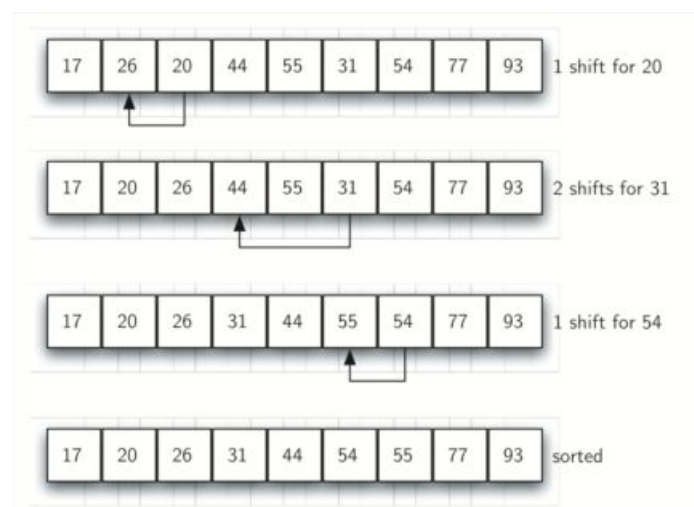


Figure 8: ShellSort: A Final Insertion Sort with Increment of 1

../\_images/shellsortD.png

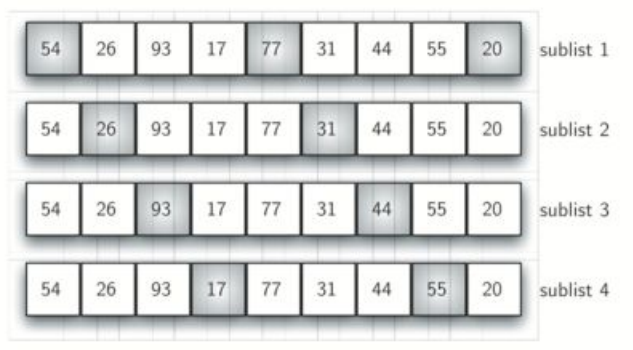


Figure 9: Initial Sublists for a Shell Sort

We said earlier that the way in which the increments are chosen is the unique feature of the shell sort. The function shown in ActiveCode 1 uses a different set of increments. In this case, we begin with  $n/2$  sublists. On the next pass,  $n/4$  sublists are sorted. Eventually, a single list is sorted with the basic insertion sort. Figure 9 shows the first sublists for our example using this increment.

The following invocation of the shellSort function shows the partially sorted lists after each increment, with the final sort being an insertion sort with an increment of one.

Run Save Load Show in Codelens

```
def shellSort(alist):
    sublistcount = len(alist)//2
    while sublistcount > 0:

        for startposition in range(sublistcount):
            gapInsertionSort(alist,startposition,sublistcount)

        print("After increments of size",sublistcount,
              "The list is",alist)

        sublistcount = sublistcount // 2

def gapInsertionSort(alist,start,gap):
    for i in range(start+gap,len(alist),gap):

        currentvalue = alist[i]
        position = i

        while position>=gap and alist[position-gap]>currentvalue:
            alist[position]=alist[position-gap]
```

```
        position = position-gap

        alist[position]=currentvalue

alist = [54,26,93,17,77,31,44,55,20]
shellSort(alist)
print(alist)
```

Shell Sort (lst\_shellSort)

Initialize Run Stop

Beginning Step Forward Step Backward End

At first glance you may think that a shell sort cannot be better than an insertion sort, since it does a complete insertion sort as the last step. It turns out, however, that this final insertion sort does not need to do very many comparisons (or shifts) since the list has been pre-sorted by earlier incremental insertion sorts, as described above. In other words, each pass produces a list that is “more sorted” than the previous one. This makes the final pass very efficient.

Although a general analysis of the shell sort is well beyond the scope of this text, we can say that it tends to fall somewhere between  $O(n)$  and  $O(n^2)$ , based on the behavior described above. For the increments shown in Listing 5, the performance is  $O(n^2)$ . By changing the increment, for example using  $2k-1$  (1, 3, 7, 15, 31, and so on), a shell sort can perform at  $O(n^3)$ .

### Self Check

Q-31: Given the following list of numbers: [5, 16, 20, 12, 3, 8, 9, 17, 19, 7] Which answer illustrates the contents of the list after all swapping is complete for a gap size of 3?

- a) [5, 3, 8, 7, 16, 19, 9, 17, 20, 12]
- b) [3, 7, 5, 8, 9, 12, 19, 16, 20, 17]
- c) [3, 5, 7, 8, 9, 12, 16, 17, 19, 20]
- d) [5, 16, 20, 3, 8, 12, 9, 17, 20, 7]

## 5.3.5. 归并排序

We now turn our attention to using a divide and conquer strategy as a way to improve the performance of sorting algorithms. The first algorithm we will study is the merge sort. Merge sort is a recursive algorithm that continually splits a list in half. If the list is empty or has one item, it is sorted by definition (the base case). If the list has more than one item, we split the list and recursively invoke a merge sort on both halves. Once the two halves are sorted, the fundamental operation, called a merge, is performed. Merging is the process of taking two smaller sorted lists and combining them together into a single, sorted, new list. Figure 10 shows our familiar example list as it is being split by mergeSort. Figure 11 shows the simple lists, now sorted, as they are merged back together.

../\_images/mergesortA.png

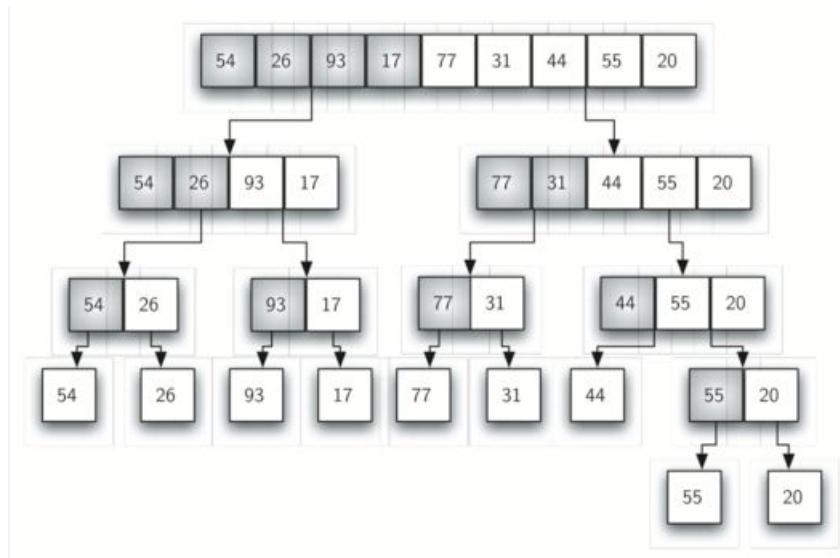


Figure 10: Splitting the List in a Merge Sort

../\_images/mergesortB.png

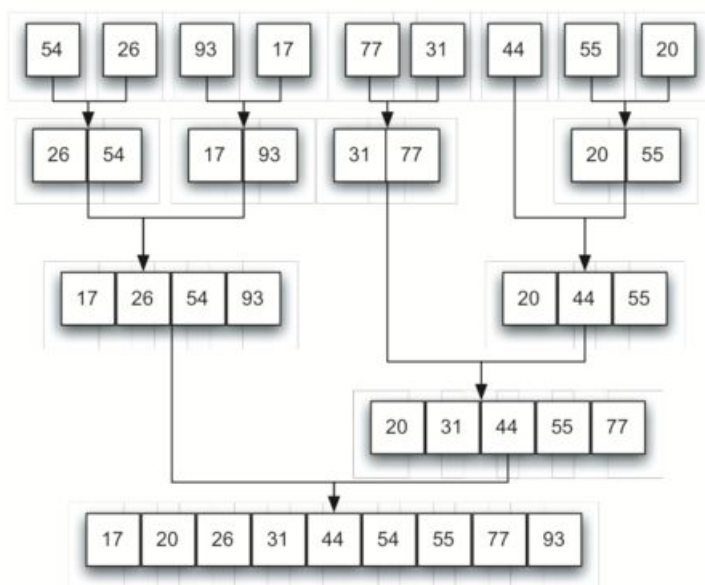


Figure 11: Lists as They Are Merged Together

The mergeSort function shown in ActiveCode 1 begins by asking the base case question. If the length of the list is less than or equal to one, then we already have a sorted list and no more processing is necessary. If, on the other hand, the length is greater than one, then we use the Python slice operation to extract the left and right halves. It is important to note that the list may not have an even number of items. That does not matter, as the lengths will differ by at most one.

Run Save Load Show in Codelens

```

def mergeSort(alist):
    print("Splitting ",alist)
    if len(alist)>1:
        mid = len(alist)//2
        lefthalf = alist[:mid]
        righthalf = alist[mid:]

        mergeSort(lefthalf)
        mergeSort(righthalf)

        i=0
        j=0
        k=0
        while i<len(lefthalf) and j<len(righthalf):
            if lefthalf[i]<righthalf[j]:
                alist[k]=lefthalf[i]
                i=i+1
            else:
                alist[k]=righthalf[j]
                j=j+1
            k=k+1

        while i<len(lefthalf):
            alist[k]=lefthalf[i]
            i=i+1
            k=k+1

        while j<len(righthalf):
            alist[k]=righthalf[j]
            j=j+1
            k=k+1
    print("Merging ",alist)

alist = [54,26,93,17,77,31,44,55,20]
mergeSort(alist)
print(alist)

```

#### Merge Sort (lst\_merge)

Once the mergeSort function is invoked on the left half and the right half (lines 8 - 9), it is assumed they are sorted. The rest of the function (lines 11 - 31) is responsible for merging the two smaller sorted lists into a larger sorted list. Notice that the merge operation places the items back into the original list (alist) one at a time by repeatedly taking the smallest item from the sorted

lists.

The mergeSort function has been augmented with a print statement (line 2) to show the contents of the list being sorted at the start of each invocation. There is also a print statement (line 32) to show the merging process. The transcript shows the result of executing the function on our example list. Note that the list with 44, 55, and 20 will not divide evenly. The first split gives [44] and the second gives [55,20]. It is easy to see how the splitting process eventually yields a list that can be immediately merged with other sorted lists.

Initialize Run Stop

Beginning Step Forward Step Backward End

In order to analyze the mergeSort function, we need to consider the two distinct processes that make up its implementation. First, the list is split into halves. We already computed (in a binary search) that we can divide a list in half  $\log n$  times where  $n$  is the length of the list. The second process is the merge. Each item in the list will eventually be processed and placed on the sorted list. So the merge operation which results in a list of size  $n$  requires  $n$  operations. The result of this analysis is that  $\log n$  splits, each of which costs  $n$  for a total of  $n \log n$  operations. A merge sort is an  $O(n \log n)$  algorithm.

Recall that the slicing operator is  $O(k)$  where  $k$  is the size of the slice. In order to guarantee that mergeSort will be  $O(n \log n)$  we will need to remove the slice operator. Again, this is possible if we simply pass the starting and ending indices along with the list when we make the recursive call. We leave this as an exercise.

It is important to notice that the mergeSort function requires extra space to hold the two halves as they are extracted with the slicing operations. This additional space can be a critical factor if the list is large and can make this sort problematic when working on large data sets.

### Self Check

Q-23: Given the following list of numbers: `[21, 1, 26, 45, 29, 28, 2, 9, 16, 49, 39, 27, 43, 34, 46, 40]` which answer illustrates the list to be sorted after 3 recursive calls to mergesort?

- a) [16, 49, 39, 27, 43, 34, 46, 40]
- b) [21,1]
- c) [21, 1, 26, 45]
- d) [21]

Check Me   Compare Me

Q-24: Given the following list of numbers: `[21, 1, 26, 45, 29, 28, 2, 9, 16, 49, 39, 27, 43, 34, 46, 40]` which answer illustrates the first two lists to be merged?

- a) [21, 1] and [26, 45]
- b) [[1, 2, 9, 21, 26, 28, 29, 45] and [16, 27, 34, 39, 40, 43, 46, 49]
- c) [21] and [1]
- d) [9] and [16]

### 5.3.6. 快速排序

The quick sort uses divide and conquer to gain the same advantages as the merge sort, while not using additional storage. As a trade-off, however, it is possible that the list may not be divided in half. When this happens, we will see that performance is diminished.

A quick sort first selects a value, which is called the pivot value. Although there are many different ways to choose the pivot value, we will simply use the first item in the list. The role of the pivot value is to assist with splitting the list. The actual position where the pivot value belongs in the final sorted list, commonly called the split point, will be used to divide the list for subsequent calls to the quick sort.

Figure 12 shows that 54 will serve as our first pivot value. Since we have looked at this example a few times already, we know that 54 will eventually end up in the position currently holding 31. The partition process will happen next. It will find the split point and at the same time move other items to the appropriate side of the list, either less than or greater than the pivot value.

../\_images/firstsplit.png

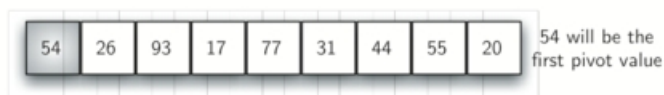


Figure 12: The First Pivot Value for a Quick Sort

Partitioning begins by locating two position markers—let's call them leftmark and rightmark—at the beginning and end of the remaining items in the list (positions 1 and 8 in Figure 13). The goal of the partition process is to move items that are on the wrong side with respect to the pivot value while also converging on the split point. Figure 13 shows this process as we locate the position of 54.

../\_images/partitionA.png



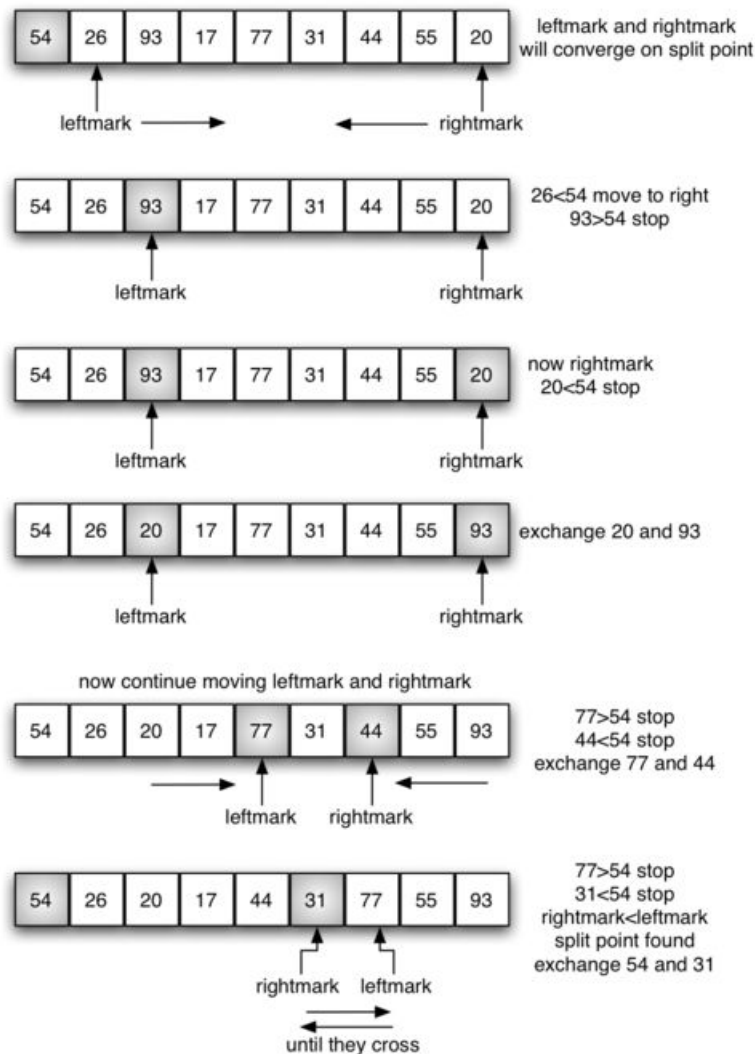


Figure 13: Finding the Split Point for 54

We begin by incrementing leftmark until we locate a value that is greater than the pivot value. We then decrement rightmark until we find a value that is less than the pivot value. At this point we have discovered two items that are out of place with respect to the eventual split point. For our example, this occurs at 93 and 20. Now we can exchange these two items and then repeat the process again.

At the point where rightmark becomes less than leftmark, we stop. The position of rightmark is now the split point. The pivot value can be exchanged with the contents of the split point and the pivot value is now in place (Figure 14). In addition, all the items to the left of the split point are less than the pivot value, and all the items to the right of the split point are greater than the pivot value. The list can now be divided at the split point and the quick sort can be invoked recursively on the two halves.

../\_images/partitionB.png

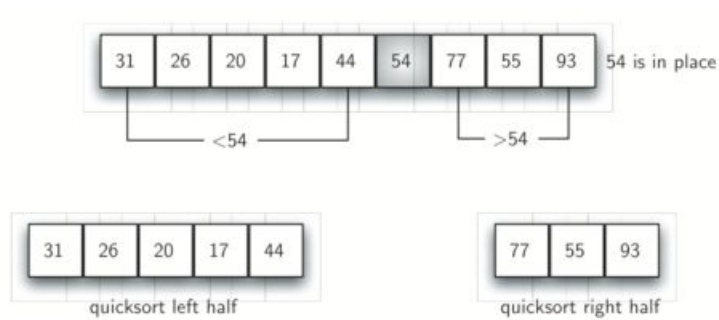


Figure 14: Completing the Partition Process to Find the Split Point for 54

The quickSort function shown in ActiveCode 1 invokes a recursive function, quickSortHelper. quickSortHelper begins with the same base case as the merge sort. If the length of the list is less than or equal to one, it is already sorted. If it is greater, then it can be partitioned and recursively sorted. The partition function implements the process described earlier.

Run Save Load Show in Codelens

```
def quickSort(alist):
    quickSortHelper(alist,0,len(alist)-1)

def quickSortHelper(alist,first,last):
    if first<last:

        splitpoint = partition(alist,first,last)

        quickSortHelper(alist,first,splitpoint-1)
        quickSortHelper(alist,splitpoint+1,last)

def partition(alist,first,last):
    pivotvalue = alist[first]

    leftmark = first+1
    rightmark = last

    done = False
    while not done:

        while leftmark <= rightmark and \
            alist[leftmark] <= pivotvalue:
            leftmark = leftmark + 1

        while alist[rightmark] >= pivotvalue and \
```

```

        rightmark >= leftmark:
            rightmark = rightmark - 1

    if rightmark < leftmark:
        done = True
    else:
        temp = alist[leftmark]
        alist[leftmark] = alist[rightmark]
        alist[rightmark] = temp

    temp = alist[first]
    alist[first] = alist[rightmark]
    alist[rightmark] = temp

    return rightmark

alist = [54,26,93,17,77,31,44,55,20]
quickSort(alist)
print(alist)

Quick Sort (1st_quick)

```

Initialize Run Stop

Beginning Step Forward Step Backward End

To analyze the quickSort function, note that for a list of length  $n$ , if the partition always occurs in the middle of the list, there will again be  $\log n$  divisions. In order to find the split point, each of the  $n$  items needs to be checked against the pivot value. The result is  $n \log n$ . In addition, there is no need for additional memory as in the merge sort process.

Unfortunately, in the worst case, the split points may not be in the middle and can be very skewed to the left or the right, leaving a very uneven division. In this case, sorting a list of  $n$  items divides into sorting a list of 0 items and a list of  $n-1$  items. Then sorting a list of  $n-1$  divides into a list of size 0 and a list of size  $n-2$ , and so on. The result is an  $O(n^2)$  sort with all of the overhead that recursion requires.

We mentioned earlier that there are different ways to choose the pivot value. In particular, we can attempt to alleviate some of the potential for an uneven division by using a technique called median of three. To choose the pivot value, we will consider the first, the middle, and the last element in the list. In our example, those are 54, 77, and 20. Now pick the median value, in our case 54, and use it for the pivot value (of course, that was the pivot value we used originally). The idea is that in the case where the first item in the list does not belong toward the middle of the list, the median of three will choose a better “middle” value. This will be particularly useful when the original list is somewhat sorted to begin with. We leave the implementation of this pivot

value selection as an exercise.

### Self Check

Q-25: Given the following list of numbers [14, 17, 13, 15, 19, 10, 3, 16, 9, 12] which answer shows the contents of the list after the second partitioning according to the quicksort algorithm?

- a) [9, 3, 10, 13, 12]
- b) [9, 3, 10, 13, 12, 14]
- c) [9, 3, 10, 13, 12, 14, 17, 16, 15, 19]
- d) [9, 3, 10, 13, 12, 14, 19, 16, 15, 17]

[Check Me](#)   [Compare Me](#)

Q-26: Given the following list of numbers [1, 20, 11, 5, 2, 9, 16, 14, 13, 19] what would be the first pivot value using the median of 3 method?

- a) 1
- b) 9
- c) 16
- d) 19

[Check Me](#)   [Compare Me](#)

Q-27: Which of the following sort algorithms are guaranteed to be  $O(n \log n)$  even in the worst case?

- a) Shell Sort
- b) Quick Sort
- c) Merge Sort
- d) Insertion Sort

## 5.4. 小结

- A sequential search is  $O(n)$  for ordered and unordered lists.
- A binary search of an ordered list is  $O(\log n)$  in the worst case.
- Hash tables can provide constant time searching.
- A bubble sort, a selection sort, and an insertion sort are  $O(n^2)$  algorithms.
- A shell sort improves on the insertion sort by sorting incremental sublists. It falls between  $O(n)$  and  $O(n^2)$ .
- A merge sort is  $O(n \log n)$ , but requires additional space for the merging process.
- A quick sort is  $O(n \log n)$ , but may degrade to  $O(n^2)$  if the split points are not near the middle of the list. It does not require additional space.

## 5.5. 关键词

binary Search	bubble Sort	chaining
clustering	collision	collision resolution

folding method	gap	hash function
hash table	hashing	insertion sort
linear probing	load factor	map
median of three	merge	merge sort
mid-square method	open addressing	partition
perfect hash function	pivot value	quadratic probing
quick sort	rehashing	selection sort
sequential search	shell sort	short bubble
slot	split point	

## 5.6. 问题讨论

1. Using the hash table performance formulas given in the chapter, compute the average number of comparisons necessary when the table is

- 10% full
- 25% full
- 50% full
- 75% full
- 90% full
- 99% full

At what point do you think the hash table is too small? Explain.

2. Modify the hash function for strings to use positional weightings.

3. We used a hash function for strings that weighted the characters by position. Devise an alternative weighting scheme. What are the biases that exist with these functions?

4. Research perfect hash functions. Using a list of names (classmates, family members, etc.), generate the hash values using the perfect hash algorithm.

5. Generate a random list of integers. Show how this list is sorted by the following algorithms:

- bubble sort
- selection sort
- insertion sort
- shell sort (you decide on the increments)
- merge sort
- quick sort (you decide on the pivot value)

6. Consider the following list of integers: [1,2,3,4,5,6,7,8,9,10]. Show how this list is sorted by the following algorithms:

- bubble sort

- selection sort
- insertion sort
- shell sort (you decide on the increments)
- merge sort
- quick sort (you decide on the pivot value)

7. Consider the following list of integers: [10,9,8,7,6,5,4,3,2,1]. Show how this list is sorted by the following algorithms:

- bubble sort
- selection sort
- insertion sort
- shell sort (you decide on the increments)
- merge sort
- quick sort (you decide on the pivot value)

8. Consider the list of characters: ['P','Y','T','H','O','N']. Show how this list is sorted using the following algorithms:

- bubble sort
- selection sort
- insertion sort
- shell sort (you decide on the increments)
- merge sort
- quick sort (you decide on the pivot value)

9. Devise alternative strategies for choosing the pivot value in quick sort. For example, pick the middle item. Re-implement the algorithm and then execute it on random data sets. Under what criteria does your new strategy perform better or worse than the strategy from this chapter?

## 5.7. 编程练习

1. Set up a random experiment to test the difference between a sequential search and a binary search on a list of integers.
2. Use the binary search functions given in the text (recursive and iterative). Generate a random, ordered list of integers and do a benchmark analysis for each one. What are your results? Can you explain them?
3. Implement the binary search using recursion without the slice operator. Recall that you will need to pass the list along with the starting and ending index values for the sublist. Generate a random, ordered list of integers and do a benchmark analysis.
4. Implement the len method (`__len__`) for the hash table Map ADT implementation.
5. Implement the in method (`__contains__`) for the hash table Map ADT implementation.
6. How can you delete items from a hash table that uses chaining for collision resolution? How about if open addressing is used? What are the special circumstances that must be handled? Implement the del method for the HashTable class.
7. In the hash table map implementation, the hash table size was chosen to be 101. If the table gets full, this needs to be increased. Re-implement the put method so that the table will

automatically resize itself when the loading factor reaches a predetermined value (you can decide the value based on your assessment of load versus performance).

8. Implement quadratic probing as a rehash technique.
9. Using a random number generator, create a list of 500 integers. Perform a benchmark analysis using some of the sorting algorithms from this chapter. What is the difference in execution speed?
10. Implement the bubble sort using simultaneous assignment.
11. A bubble sort can be modified to “bubble” in both directions. The first pass moves “up” the list, and the second pass moves “down.” This alternating pattern continues until no more passes are necessary. Implement this variation and describe under what circumstances it might be appropriate.
12. Implement the selection sort using simultaneous assignment.
13. Perform a benchmark analysis for a shell sort, using different increment sets on the same list.
14. Implement the mergeSort function without using the slice operator.
15. One way to improve the quick sort is to use an insertion sort on lists that have a small length (call it the “partition limit” ). Why does this make sense? Re-implement the quick sort and use it to sort a random list of integers. Perform an analysis using different list sizes for the partition limit.
16. Implement the median-of-three method for selecting a pivot value as a modification to quickSort. Run an experiment to compare the two techniques.

## 6. 树和树算法

### 6.1. 目标

- 了解树的数据结构及使用；
- 树用来实现映射 `map` 数据结构；
- 采用 `List` 来实现树；
- 采用类和引用来实现树；
- 树实现为递归数据结构；
- 采用堆 `heap` 实现优先队列。



## **6.2. 树的例子**

## **6.3. 术语表与定义**

### **6.3.1. 术语表**

### **6.3.2. 定义**

## **6.4. 实现**

### **6.4.1. “列表的列表”表示树**

### **6.4.2. 节点和引用**

## **6.5. 二叉堆 Binary Heap 实现的优先队列**

### **6.5.1. 二叉堆操作**

### **6.5.2. 二叉堆实现**

## **6.6. 二叉树应用**

## **6.7. 树遍历**

## **6.8. 二叉搜索树**

### **6.8.1. 搜索树操作**

### **6.8.2. 搜索树实现**

### **6.8.3. 搜索树分析**

## **6.9. 小结**

## **6.10. 关键词**

## **6.11. 问题讨论**

## **6.12. 编程练习**

# **7. 图和图算法**

## **7.1. 目标**

## **7.2. 词汇表及定义**

## **7.3. 图抽象数据类型**

## **7.4. 邻接矩阵**

## **7.5. 邻接表**

## **7.6. 实现**

## **7.7. Word Ladder 词梯问题**

### **7.7.1. 建立 Word Ladder 图**

### **7.7.2. 实现广度优先搜索**

### **7.7.3. 广度优先搜索分析**

## **7.8. 骑士周游问题**

### **7.8.1. 建立骑士周游图**

### **7.8.2. 实现骑士周游**

### **7.8.3. 骑士周游分析**

### **7.8.4. 通用深度优先搜索**

### **7.8.5. 深度优先分析**

## **7.9. 拓扑排序**

## **7.10. 强连通分支**

## **7.11. 最短路径问题**

### **7.11.1. Dijkstra 算法**

### **7.11.2. Dijkstra 算法分析**

## **7.12. Prim 最小生成树算法**

## **7.13. 小结**

## **7.14. 关键词**

## **7.15. 问题讨论**

## **7.16. 编程练习**

