

---

# How to Think Like a Computer Scientist: Learning with Python 3 Documentation

***Release 3rd Edition***

Follow me on [LinkedIn](#) for more:  
**Steve Nouri**

<https://www.linkedin.com/in/stevenouri/>

**Peter Wentworth, Jeffrey Elkner,  
Allen B. Downey and Chris Meyers**

Feb 26, **2019**

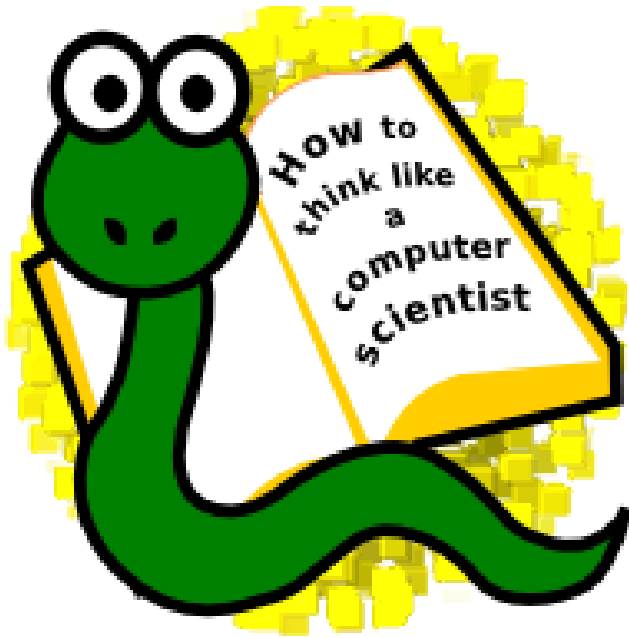
---

## Contents

---

<b>1</b>	<b>The way of the program</b>	<b>3</b>
<b>2</b>	<b>Variables, expressions and statements</b>	<b>11</b>
<b>3</b>	<b>Program Flow</b>	<b>23</b>
<b>4</b>	<b>Functions</b>	<b>63</b>
<b>5</b>	<b>Data Types</b>	<b>91</b>
<b>6</b>	<b>Numpy</b>	<b>133</b>
<b>7</b>	<b>Files</b>	<b>139</b>
<b>8</b>	<b>Modules</b>	<b>145</b>
<b>9</b>	<b>More datatypes</b>	<b>157</b>
<b>10</b>	<b>Recursion</b>	<b>161</b>
<b>11</b>	<b>Classes and Objects</b>	<b>175</b>
<b>12</b>	<b>Exceptions</b>	<b>213</b>
<b>13</b>	<b>Fitting</b>	<b>217</b>
<b>14</b>	<b>PyGame</b>	<b>223</b>
<b>15</b>	<b>Copyright Notice</b>	<b>243</b>
<b>16</b>	<b>Contributions</b>	<b>245</b>
<b>A</b>	<b>Modules</b>	<b>249</b>
<b>B</b>	<b>More datatypes</b>	<b>261</b>
<b>C</b>	<b>Recursion</b>	<b>265</b>
<b>D</b>	<b>Classes and Objects</b>	<b>279</b>

<b>E</b>	<b>Exceptions</b>	<b>317</b>
<b>F</b>	<b>Fitting</b>	<b>321</b>
<b>G</b>	<b>PyGame</b>	<b>327</b>
<b>H</b>	<b>Plotting data with matplotlib</b>	<b>347</b>



3rd Edition (Using Python 3.x)

by Jeffrey Elkner, Peter Wentworth, Allen B. Downey, and Chris Meyers

illustrated by Dario Mitchell

- [Copyright Notice](#)
- [Contributor List](#)
- [Chapter 1](#) *The way of the program*
- [Chapter 2](#) *Variables, expressions, and statements*
- [Chapter 3](#) *Program Flow*
- [Chapter 4](#) *Functions*
- [Chapter 5](#) *Datatypes*
- [Chapter 6](#) *Numpy*
- [Chapter 7](#) *File I/O*
- [Appendix A](#) *Writing Your Own Modules*
- [Appendix B](#) *More datatypes*
- [Appendix C](#) *Recursion*
- [Appendix D](#) *Object Oriented Programming*
- [Appendix E](#) *Exceptions*
- [Appendix F](#) *Fitting and Scientific Data Handling*
- [Appendix G](#) *PyGame*
- [Appendix H](#) *Plotting with matplotlib*
- [GNU Free Document License](#)



# CHAPTER 1

---

## The way of the program

---

The goal of this book is to teach you to think like a computer scientist. This way of thinking combines some of the best features of mathematics, engineering, and natural science. Like mathematicians, computer scientists use formal languages to denote ideas (specifically computations). Like engineers, they design things, assembling components into systems and evaluating tradeoffs among alternatives. Like scientists, they observe the behavior of complex systems, form hypotheses, and test predictions.

The single most important skill for a computer scientist is **problem solving**. Problem solving means the ability to formulate problems, think creatively about solutions, and express a solution clearly and accurately. As it turns out, the process of learning to program is an excellent opportunity to practice problem-solving skills. That's why this chapter is called, The way of the program.

On one level, you will be learning to program, a useful skill by itself. On another level, you will use programming as a means to an end. As we go along, that end will become clearer.

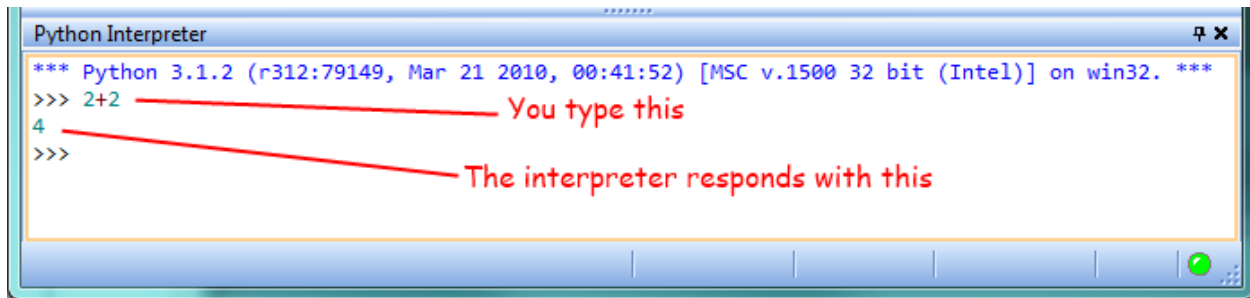
### 1.1 The Python programming language

The programming language you will be learning is Python. Python is an example of a **high-level language**; other high-level languages you might have heard of are C++, PHP, Pascal, C#, and Java.

As you might infer from the name high-level language, there are also **low-level languages**, sometimes referred to as machine languages or assembly languages. Loosely speaking, computers can only execute programs written in low-level languages. Thus, programs written in a high-level language have to be translated into something more suitable before they can run.

Almost all programs are written in high-level languages because of their advantages. It is much easier to program in a high-level language so programs take less time to write, they are shorter and easier to read, and they are more likely to be correct. Second, high-level languages are **portable**, meaning that they can run on different kinds of computers with few or no modifications.

The engine that translates and runs Python is called the **Python Interpreter**: There are two ways to use it: *immediate mode* and *script mode*. In immediate mode, you type Python expressions into the Python Interpreter window, and the interpreter immediately shows the result:



The `>>>` is called the **Python prompt**. The interpreter uses the prompt to indicate that it is ready for instructions. We typed `2 + 2`, and the interpreter evaluated our expression, and replied `4`, and on the next line it gave a new prompt, indicating that it is ready for more input.

Alternatively, you can write a program in a file and use the interpreter to execute the contents of the file. Such a file is called a **script**. Scripts have the advantage that they can be saved to disk, printed, and so on.

Working directly in the interpreter is convenient for testing short bits of code because you get immediate feedback. Think of it as scratch paper used to help you work out problems. Anything longer than a few lines should be put into a script.

## 1.2 What is a program?

A **program** is a sequence of instructions that specifies how to perform a computation. The computation might be something mathematical, such as solving a system of equations or finding the roots of a polynomial, but it can also be a symbolic computation, such as searching and replacing text in a document or (strangely enough) compiling a program.

The details look different in different languages, but a few basic instructions appear in just about every language:

**input** Get data from the keyboard, a file, or some other device such as a sensor.

**output** Display data on the screen or send data to a file or other device such as a motor.

**math** Perform basic mathematical operations like addition and multiplication.

**conditional execution** Check for certain conditions and execute the appropriate sequence of statements.

**repetition** Perform some action repeatedly, usually with some variation.

Believe it or not, that's pretty much all there is to it. Every program you've ever used, no matter how complicated, is made up of instructions that look more or less like these. Thus, we can describe programming as the process of breaking a large, complex task into smaller and smaller subtasks until the subtasks are simple enough to be performed with sequences of these basic instructions.

That may be a little vague, but we will come back to this topic later when we talk about **algorithms**.

## 1.3 What is debugging?

Programming is a complex process, and because it is done by human beings, it often leads to errors. Programming errors are called **bugs** and the process of tracking them down and correcting them is called **debugging**. Use of the term *bug* to describe small engineering difficulties dates back to at least 1889, when Thomas Edison had a bug with his phonograph.

Three kinds of errors can occur in a program: **syntax errors**, **runtime errors**, and **semantic errors**. It is useful to distinguish between them in order to track them down more quickly.

## 1.4 Syntax errors

Python can only execute a program if the program is syntactically correct; otherwise, the process fails and returns an error message. **Syntax** refers to the structure of a program and the rules about that structure. For example, in English, a sentence must begin with a capital letter and end with a period. This sentence contains a **syntax error**. So does this one

For most readers, a few syntax errors are not a significant problem, which is why we can read the poetry of E. E. Cummings without problems. Python is not so forgiving. If there is a single syntax error anywhere in your program, Python will display an error message and quit, and you will not be able to run your program. During the first few weeks of your programming career, you will probably spend a lot of time tracking down syntax errors. As you gain experience, though, you will make fewer errors and find them faster.

## 1.5 Runtime errors

The second type of error is a runtime error, so called because the error does not appear until you run the program. These errors are also called **exceptions** because they usually indicate that something exceptional (and bad) has happened.

Runtime errors are rare in the simple programs you will see in the first few chapters, so it might be a while before you encounter one.

## 1.6 Semantic errors

The third type of error is the **semantic error**. If there is a semantic error in your program, it will run successfully, in the sense that the computer will not generate any error messages, but it will not do the right thing. It will do something else. Specifically, it will do what you told it to do.

The problem is that the program you wrote is not the program you wanted to write. The meaning of the program (its semantics) is wrong. Identifying semantic errors can be tricky because it requires you to work backward by looking at the output of the program and trying to figure out what it is doing.

## 1.7 Experimental debugging

One of the most important skills you will acquire is debugging. Although it can be frustrating, debugging is one of the most intellectually rich, challenging, and interesting parts of programming.

In some ways, debugging is like detective work. You are confronted with clues, and you have to infer the processes and events that led to the results you see.

Debugging is also like an experimental science. Once you have an idea what is going wrong, you modify your program and try again. If your hypothesis was correct, then you can predict the result of the modification, and you take a step closer to a working program. If your hypothesis was wrong, you have to come up with a new one. As Sherlock Holmes pointed out, When you have eliminated the impossible, whatever remains, however improbable, must be the truth. (A. Conan Doyle, *The Sign of Four*)

For some people, programming and debugging are the same thing. That is, programming is the process of gradually debugging a program until it does what you want. The idea is that you should start with a program that does *something* and make small modifications, debugging them as you go, so that you always have a working program.

For example, Linux is an operating system kernel that contains millions of lines of code, but it started out as a simple program Linus Torvalds used to explore the Intel 80386 chip. According to Larry Greenfield, one of Linus's earlier



projects was a program that would switch between displaying AAAA and BBBB. This later evolved to Linux (*The Linux Users' Guide Beta Version 1*).

Later chapters will make more suggestions about debugging and other programming practices.

## 1.8 Formal and natural languages

**Natural languages** are the languages that people speak, such as English, Spanish, and French. They were not designed by people (although people try to impose some order on them); they evolved naturally.

**Formal languages** are languages that are designed by people for specific applications. For example, the notation that mathematicians use is a formal language that is particularly good at denoting relationships among numbers and symbols. Chemists use a formal language to represent the chemical structure of molecules. And most importantly:

*Programming languages are formal languages that have been designed to express computations.*

Formal languages tend to have strict rules about syntax. For example,  $3+3=6$  is a syntactically correct mathematical statement, but  $3=+6\$$  is not.  $\text{H}_2\text{O}$  is a syntactically correct chemical name, but  ${}_2\text{Zz}$  is not.

Syntax rules come in two flavors, pertaining to **tokens** and structure. Tokens are the basic elements of the language, such as words, numbers, parentheses, commas, and so on. In Python, a statement like `print("Happy New Year for ", 2013)` has 6 tokens: a function name, an open parenthesis (round bracket), a string, a comma, a number, and a close parenthesis.

It is possible to make errors in the way one constructs tokens. One of the problems with  $3=+6\$$  is that  $\$$  is not a legal token in mathematics (at least as far as we know). Similarly,  ${}_2\text{Zz}$  is not a legal token in chemistry notation because there is no element with the abbreviation Zz.

The second type of syntax rule pertains to the **structure** of a statement—that is, the way the tokens are arranged. The statement  $3=+6\$$  is structurally illegal because you can't place a plus sign immediately after an equal sign. Similarly, molecular formulas have to have subscripts after the element name, not before. And in our Python example, if we omitted the comma, or if we changed the two parentheses around to say `print)"Happy New Year for ", 2013 (` our statement would still have six legal and valid tokens, but the structure is illegal.

When you read a sentence in English or a statement in a formal language, you have to figure out what the structure of the sentence is (although in a natural language you do this subconsciously). This process is called **parsing**.

For example, when you hear the sentence, “The other shoe fell”, you understand that the other shoe is the subject and fell is the verb. Once you have parsed a sentence, you can figure out what it means, or the **semantics** of the sentence. Assuming that you know what a shoe is and what it means to fall, you will understand the general implication of this sentence.

Although formal and natural languages have many features in common — tokens, structure, syntax, and semantics — there are many differences:

**ambiguity** Natural languages are full of ambiguity, which people deal with by using contextual clues and other information. Formal languages are designed to be nearly or completely unambiguous, which means that any statement has exactly one meaning, regardless of context.

**redundancy** In order to make up for ambiguity and reduce misunderstandings, natural languages employ lots of redundancy. As a result, they are often verbose. Formal languages are less redundant and more concise.

**literalness** Formal languages mean exactly what they say. On the other hand, natural languages are full of idiom and metaphor. If someone says, “The other shoe fell”, there is probably no shoe and nothing falling. You'll need to find the original joke to understand the idiomatic meaning of the other shoe falling. *Yahoo! Answers* thinks it knows!

People who grow up speaking a natural language—everyone—often have a hard time adjusting to formal languages. In some ways, the difference between formal and natural language is like the difference between poetry and prose, but more so:

**poetry** Words are used for their sounds as well as for their meaning, and the whole poem together creates an effect or emotional response. Ambiguity is not only common but often deliberate.

**prose** The literal meaning of words is more important, and the structure contributes more meaning. Prose is more amenable to analysis than poetry but still often ambiguous.

**program** The meaning of a computer program is unambiguous and literal, and can be understood entirely by analysis of the tokens and structure.

Here are some suggestions for reading programs (and other formal languages). First, remember that formal languages are much more dense than natural languages, so it takes longer to read them. Also, the structure is very important, so it is usually not a good idea to read from top to bottom, left to right. Instead, learn to parse the program in your head, identifying the tokens and interpreting the structure. Finally, the details matter. Little things like spelling errors and bad punctuation, which you can get away with in natural languages, can make a big difference in a formal language.

## 1.9 The first program

Traditionally, the first program written in a new language is called *Hello, World!* because all it does is display the words, Hello, World! In Python, the script looks like this: (For scripts, we'll show line numbers to the left of the Python statements.)

```
1 print("Hello, World!")
```

This is an example of using the **print function**, which doesn't actually print anything on paper. It displays a value on the screen. In this case, the result shown is

```
1 Hello, World!
```

The quotation marks in the program mark the beginning and end of the value; they don't appear in the result.

Some people judge the quality of a programming language by the simplicity of the Hello, World! program. By this standard, Python does about as well as possible.

## 1.10 Comments

As programs get bigger and more complicated, they get more difficult to read. Formal languages are dense, and it is often difficult to look at a piece of code and figure out what it is doing, or why.

For this reason, it is a good idea to add notes to your programs to explain in natural language what the program is doing.

A **comment** in a computer program is text that is intended only for the human reader — it is completely ignored by the interpreter.

In Python, the `#` token starts a comment. The rest of the line is ignored. Here is a new version of *Hello, World!*.

```
1 #-----
2 # This demo program shows off how elegant Python is!
3 # Written by Joe Soap, December 2010.
4 # Anyone may freely copy or modify this program.
5 #-----
```

(continues on next page)

(continued from previous page)

```
6
7 print("Hello, World!")      # Isn't this easy!
```

You'll also notice that we've left a blank line in the program. Blank lines are also ignored by the interpreter, but comments and blank lines can make your programs much easier for humans to parse. Use them liberally!

## 1.11 Glossary

**algorithm** A set of specific steps for solving a category of problems.

**bug** An error in a program.

**comment** Information in a program that is meant for other programmers (or anyone reading the source code) and has no effect on the execution of the program.

**debugging** The process of finding and removing any of the three kinds of programming errors.

**exception** Another name for a runtime error.

**formal language** Any one of the languages that people have designed for specific purposes, such as representing mathematical ideas or computer programs; all programming languages are formal languages.

**high-level language** A programming language like Python that is designed to be easy for humans to read and write.

**immediate mode** A style of using Python where we type expressions at the command prompt, and the results are shown immediately. Contrast with **script**, and see the entry under **Python shell**.

**interpreter** The engine that executes your Python scripts or expressions.

**low-level language** A programming language that is designed to be easy for a computer to execute; also called machine language or assembly language.

**natural language** Any one of the languages that people speak that evolved naturally.

**object code** The output of the compiler after it translates the program.

**parse** To examine a program and analyze the syntactic structure.

**portability** A property of a program that can run on more than one kind of computer.

**print function** A function used in a program or script that causes the Python interpreter to display a value on its output device.

**problem solving** The process of formulating a problem, finding a solution, and expressing the solution.

**program** a sequence of instructions that specifies to a computer actions and computations to be performed.

**Python shell** An interactive user interface to the Python interpreter. The user of a Python shell types commands at the prompt (`>>>`), and presses the return key to send these commands immediately to the interpreter for processing. The word *shell* comes from Unix.

**runtime error** An error that does not occur until the program has started to execute but that prevents the program from continuing.

**script** A program stored in a file (usually one that will be interpreted).

**semantic error** An error in a program that makes it do something other than what the programmer intended.

**semantics** The meaning of a program.

**source code** A program in a high-level language before being compiled.

**syntax** The structure of a program.

**syntax error** An error in a program that makes it impossible to parse — and therefore impossible to interpret.

**token** One of the basic elements of the syntactic structure of a program, analogous to a word in a natural language.

## 1.12 Exercises

1. Write an English (or Dutch!) sentence with understandable semantics but incorrect syntax. Write another English (or Dutch!) sentence which has correct syntax but has semantic errors.
2. Using the Python interpreter, type `1 + 2` and then hit return. Python *evaluates* this *expression*, displays the result, and then shows another prompt. `*` is the *multiplication operator*, and `**` is the *exponentiation operator*. Experiment by entering different expressions and recording what is displayed by the Python interpreter.
3. Type `1 2` and then hit return. Python tries to evaluate the expression, but it can't because the expression is not syntactically legal. Instead, it shows the error message:

```
File "<interactive input>", line 1
 1 2
   ^
SyntaxError: invalid syntax
```

In many cases, Python indicates where the syntax error occurred, but it is not always right, and it doesn't give you much information about what is wrong.

So, for the most part, the burden is on you to learn the syntax rules.

In this case, Python is complaining because there is no operator between the numbers.

See if you can find a few more examples of things that will produce error messages when you enter them at the Python prompt. Write down what you enter at the prompt and the last line of the error message that Python reports back to you.

4. Type `print("hello")`. Python executes this, which has the effect of printing the letters h-e-l-l-o. Notice that the quotation marks that you used to enclose the string are not part of the output. Now type `"hello"` and describe your result. Make notes of when you see the quotation marks and when you don't.
5. Type `cheese` without the quotation marks. The output will look something like this:

```
Traceback (most recent call last):
  File "<interactive input>", line 1, in ?
NameError: name 'cheese' is not defined
```

This is a run-time error; specifically, it is a `NameError`, and even more specifically, it is an error because the name *cheese* is not defined. If you don't know what that means yet, you will soon.

6. Type `6 + 4 * 9` at the Python prompt and hit enter. Record what happens.

Now create a Python script with the following contents:

```
1 6 + 4 * 9
```

What happens when you run this script? Now change the script contents to:

```
1 print(6 + 4 * 9)
```

and run it again.

What happened this time?

Whenever an *expression* is typed at the Python prompt, it is evaluated and the result is *automatically* shown on the line below. (Like on your calculator, if you type this expression you'll get the result 42.)

A script is different, however. Evaluations of expressions are not automatically displayed, so it is necessary to use the **print** function to make the answer show up.

It is hardly ever necessary to use the print function in immediate mode at the command prompt.

---

## Variables, expressions and statements

---

### 2.1 Values and data types

A **value** is one of the fundamental things — like a letter or a number — that a program manipulates. The values we have seen so far are 4 (the result when we added  $2 + 2$ ), and "Hello, World!".

These values are classified into different **classes**, or **data types**: 4 is an *integer*, and "Hello, World!" is a *string*, so-called because it contains a string of letters. You (and the interpreter) can identify strings because they are enclosed in quotation marks.

If you are not sure what class a value falls into, Python has a function called **type** which can tell you.

```
>>> type("Hello, World!")
<class 'str'>
>>> type(17)
<class 'int'>
```

Not surprisingly, strings belong to the class **str** and integers belong to the class **int**. Less obviously, numbers with a decimal point belong to a class called **float**, because these numbers are represented in a format called *floating-point*. At this stage, you can treat the words *class* and *type* interchangeably. We'll come back to a deeper understanding of what a class is in later chapters.

```
>>> type(3.2)
<class 'float'>
```

What about values like "17" and "3.2"? They look like numbers, but they are in quotation marks like strings.

```
>>> type("17")
<class 'str'>
>>> type("3.2")
<class 'str'>
```

They're strings!

Strings in Python can be enclosed in either single quotes (') or double quotes ("), or three of each (''' or """)

```
>>> type('This is a string.')
<class 'str'>
>>> type("And so is this.")
<class 'str'>
>>> type("""and this.""")
<class 'str'>
>>> type(''and even this...'')
<class 'str'>
```

Double quoted strings can contain single quotes inside them, as in "Bruce's beard", and single quoted strings can have double quotes inside them, as in 'The knights who say "Ni!"'.

Strings enclosed with three occurrences of either quote symbol are called triple quoted strings. They can contain either single or double quotes:

```
>>> print('"'Oh no", she exclaimed, "Ben's bike is broken!'"')
" Oh no", she exclaimed, "Ben's bike is broken!"
>>>
```

Triple quoted strings can even span multiple lines:

```
>>> message = """This message will
... span several
... lines."""
>>> print(message)
This message will
span several
lines.
>>>
```

Python doesn't care whether you use single or double quotes or the three-of-a-kind quotes to surround your strings: once it has parsed the text of your program or command, the way it stores the value is identical in all cases, and the surrounding quotes are not part of the value. But when the interpreter wants to display a string, it has to decide which quotes to use to make it look like a string.

```
>>> 'This is a string.'
'This is a string.'
>>> """And so is this."""
'And so is this.'
```

So the Python language designers usually chose to surround their strings by single quotes. What do think would happen if the string already contained single quotes?

When you type a large integer, you might be tempted to use commas between groups of three digits, as in 42,000. The same goes for entering Dutch-style floating point numbers using a comma instead of a decimal dot. This is not a legal integer in Python, but it does mean something else, which is legal:

```
>>> 42000
42000
>>> 42,000
(42, 0)
```

Well, that's not what we expected at all! Because of the comma, Python chose to treat this as a *pair* of values. We'll come back to learn about pairs later. But, for the moment, remember not to put commas or spaces in your integers, no matter how big they are. Also revisit what we said in the previous chapter: formal languages are strict, the notation is concise, and even the smallest change might mean something quite different from what you intended.

## 2.2 Variables

One of the most powerful features of a programming language is the ability to manipulate **variables**. A variable is a name that refers to a value.

The **assignment statement** gives a value to a variable:

```
>>> message = "What's up, Doc?"
>>> n = 17
>>> pi = 3.14159
```

This example makes three assignments. The first assigns the string value "What's up, Doc?" to a variable named `message`. The second gives the integer 17 to `n`, and the third assigns the floating-point number 3.14159 to a variable called `pi`.

The **assignment token**, `=`, should not be confused with *equals*, which uses the token `==`. The assignment statement binds a *name*, on the left-hand side of the operator, to a *value*, on the right-hand side. Basically, an assignment is an order, and the equals operator can be read as a question mark. This is why you will get an error if you enter:

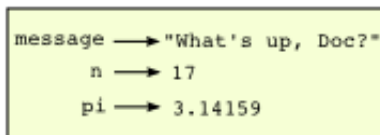
```
>>> 17 = n
File "<interactive input>", line 1
SyntaxError: can't assign to literal
```

---

**Tip:** When reading or writing code, say to yourself “`n` is assigned 17” or “`n` gets the value 17”. Don’t say “`n` equals 17”.

---

A common way to represent variables on paper is to write the name with an arrow pointing to the variable’s value. This kind of figure is called a **state snapshot** because it shows what state each of the variables is in at a particular instant in time. (Think of it as the variable’s state of mind). Some editors and programming environments do this for you, and allow you to click through the state of the program saving you some paper. This diagram shows the result of executing the assignment statements:



If you ask the interpreter to evaluate a variable, it will produce the value that is currently linked to the variable:

```
>>> message
'What's up, Doc?'
>>> n
17
>>> pi
3.14159
```

We use variables in a program to “remember” things, perhaps the current score at the football game. But variables are *variable*. This means they can change over time, just like the scoreboard at a football game. You can assign a value to a variable, and later assign a different value to the same variable. (*This is different from maths. In maths, if you give ‘ $x$ ’ the value 3, it cannot change to link to a different value half-way through your calculations!*)

```
>>> day = "Thursday"
>>> day
'Thursday'
```

(continues on next page)



(continued from previous page)

```
>>> day = "Friday"
>>> day
'Friday'
>>> day = 21
>>> day
21
```

You'll notice we changed the value of `day` three times, and on the third assignment we even made it refer to a value that was of a different type.

A great deal of programming is about having the computer remember things, e.g. *The number of missed calls on your phone*, and then arranging to update or change the variable when you miss another call.

## 2.3 Variable names and keywords

**Variable names** can be arbitrarily long. They can contain both letters and digits, but they have to begin with a letter or an underscore. Although it is legal to use uppercase letters, by convention we don't. If you do, remember that case matters. `Bruce` and `bruce` are different variables.

The underscore character (`_`) can appear in a name. It is often used in names with multiple words, such as `my_name` or `price_of_tea_in_china`.

There are some situations in which names beginning with an underscore have special meaning, so a safe rule for beginners is to start all names with a letter.

If you give a variable an illegal name, you get a syntax error:

```
>>> 76trombones = "big parade"
SyntaxError: invalid syntax
>>> more$ = 1000000
SyntaxError: invalid syntax
>>> class = "Computer Science 101"
SyntaxError: invalid syntax
```

`76trombones` is illegal because it does not begin with a letter. `more$` is illegal because it contains an illegal character, the dollar sign. But what's wrong with `class`?

It turns out that `class` is one of the Python **keywords**. Keywords define the language's syntax rules and structure, and they cannot be used as variable names.

Python has thirty-something keywords (and every now and again improvements to Python introduce or eliminate one or two):

<code>and</code>	<code>as</code>	<code>assert</code>	<code>break</code>	<code>class</code>	<code>continue</code>
<code>def</code>	<code>del</code>	<code>elif</code>	<code>else</code>	<code>except</code>	<code>exec</code>
<code>finally</code>	<code>for</code>	<code>from</code>	<code>global</code>	<code>if</code>	<code>import</code>
<code>in</code>	<code>is</code>	<code>lambda</code>	<code>nonlocal</code>	<code>not</code>	<code>or</code>
<code>pass</code>	<code>raise</code>	<code>return</code>	<code>try</code>	<code>while</code>	<code>with</code>
<code>yield</code>	<code>True</code>	<code>False</code>	<code>None</code>		

You might want to keep this list handy. If the interpreter complains about one of your variable names and you don't know why, see if it is on this list.

Programmers generally choose names for their variables that are meaningful to the human readers of the program — they help the programmer document, or remember, what the variable is used for.

**Caution:** Beginners sometimes confuse “meaningful to the human readers” with “meaningful to the computer”. So they’ll wrongly think that because they’ve called some variable `average` or `pi`, it will somehow magically calculate an average, or magically know that the variable `pi` should have a value like 3.14159. No! The computer doesn’t understand what you intend the variable to mean.

So you’ll find some instructors who deliberately don’t choose meaningful names when they teach beginners — not because we don’t think it is a good habit, but because we’re trying to reinforce the message that you — the programmer — must write the program code to calculate the average, and you must write an assignment statement to give the variable `pi` the value you want it to have.

```
e = 3.1415
ray = 10
size = e * ray ** 2
```

```
pi = 3.1415
radius = 10
area = pi * radius ** 2
```

The above two snippets do exactly the same thing, but the bottom one uses the right kind of variable names. For the computer there is no difference at all, but for a human, using the names and letters that are part of the conventional way of writing things make all the difference in the world. Using `e` instead of `pi` completely confuses people, while computers will just perform the calculation!

## 2.4 Statements

A **statement** is an instruction that the Python interpreter can execute. We have only seen the assignment statement so far. Some other kinds of statements that we’ll see shortly are `while` statements, `for` statements, `if` statements, and `import` statements. (There are other kinds too!)

When you type a statement on the command line, Python executes it. Statements don’t produce any result.

## 2.5 Evaluating expressions

An **expression** is a combination of values, variables, operators, and calls to functions. If you type an expression at the Python prompt, the interpreter **evaluates** it and displays the result:

```
>>> 1 + 1
2
>>> len("hello")
5
```

In this example `len` is a built-in Python function that returns the number of characters in a string. We’ve previously seen the `print` and the `type` functions, so this is our third example of a function!

The *evaluation of an expression* produces a value, which is why expressions can appear on the right hand side of assignment statements. A value all by itself is a simple expression, and so is a variable.

```
>>> 17
17
>>> y = 3.14
>>> x = len("hello")
>>> x
```

(continues on next page)

(continued from previous page)

```
5
>>> y
3.14
```

## 2.6 Operators and operands

**Operators** are special tokens that represent computations like addition, multiplication and division. The values the operator uses are called **operands**.

The following are all legal Python expressions whose meaning is more or less clear:

```
20+32    hour-1    hour*60+minute    minute/60    5**2    (5+9)*(15-7)
```

The tokens `+`, `-`, and `*`, and the use of parenthesis for grouping, mean in Python what they mean in mathematics. The asterisk (`*`) is the token for multiplication, and `**` is the token for exponentiation.

```
>>> 2 ** 3
8
>>> 3 ** 2
9
```

When a variable name appears in the place of an operand, it is replaced with its value before the operation is performed.

Addition, subtraction, multiplication, and exponentiation all do what you expect.

Example: so let us convert 645 minutes into hours:

```
>>> minutes = 645
>>> hours = minutes / 60
>>> hours
10.75
```

Oops! In Python 3, the division operator `/` always yields a floating point result. What we might have wanted to know was how many *whole* hours there are, and how many minutes remain. Python gives us two different flavors of the division operator. The second, called **floor division** uses the token `//`. Its result is always a whole number — and if it has to adjust the number it always moves it to the left on the number line. So `6//4` yields `1`, but `-6//4` might surprise you!

```
>>> 7 / 4
1.75
>>> 7 // 4
1
>>> minutes = 645
>>> hours = minutes // 60
>>> hours
10
```

Take care that you choose the correct flavor of the division operator. If you're working with expressions where you need floating point values, use the division operator that does the division accurately.

## 2.7 Type converter functions

Here we'll look at three more Python functions, `int`, `float` and `str`, which will (attempt to) convert their arguments into types `int`, `float` and `str` respectively. We call these **type converter** functions.

The `int` function can take a floating point number or a string, and turn it into an `int`. For floating point numbers, it *discards* the decimal portion of the number — a process we call *truncation towards zero* on the number line. Let us see this in action:

```
>>> int(3.14)
3
>>> int(3.9999)           # This doesn't round to the closest int!
3
>>> int(3.0)
3
>>> int(-3.999)           # Note that the result is closer to zero
-3
>>> int(minutes / 60)
10
>>> int("2345")           # Parse a string to produce an int
2345
>>> int(17)               # It even works if arg is already an int
17
>>> int("23 bottles")
```

This last case doesn't look like a number — what do we expect?

```
Traceback (most recent call last):
  File "<interactive input>", line 1, in <module>
ValueError: invalid literal for int() with base 10: '23 bottles'
```

The type converter `float` can turn an integer, a float, or a syntactically legal string into a float:

```
>>> float(17)
17.0
>>> float("123.45")
123.45
```

The type converter `str` turns its argument into a string:

```
>>> str(17)
'17'
>>> str(123.45)
'123.45'
```

## 2.8 Order of operations

When more than one operator appears in an expression, the order of evaluation depends on the **rules of precedence**. Python follows the same precedence rules for its mathematical operators that mathematics does. The acronym PEMDAS is a useful way to remember the order of operations:

1. **Parentheses** have the highest precedence and can be used to force an expression to evaluate in the order you want. Since expressions in parentheses are evaluated first,  $2 * (3-1)$  is 4, and  $(1+1) ** (5-2)$  is 8. You can also use parentheses to make an expression easier to read, as in  $(\text{minute} * 100) / 60$ , even though it doesn't change the result.

2. Exponentiation has the next highest precedence, so  $2**1+1$  is 3 and not 4, and  $3*1**3$  is 3 and not 27.
3. Multiplication and both Division operators have the same precedence, which is higher than Addition and Subtraction, which also have the same precedence. So  $2*3-1$  yields 5 rather than 4, and  $5-2*2$  is 1, not 6.
4. Operators with the *same* precedence are evaluated from left-to-right. In algebra we say they are *left-associative*. So in the expression  $6-3+2$ , the subtraction happens first, yielding 3. We then add 2 to get the result 5. If the operations had been evaluated from right to left, the result would have been  $6-(3+2)$ , which is 1. (The acronym PEDMAS could mislead you to thinking that division has higher precedence than multiplication, and addition is done ahead of subtraction - don't be misled. Subtraction and addition are at the same precedence, and the left-to-right rule applies.)

Due to some historical quirk, an exception to the left-to-right left-associative rule is the exponentiation operator `**`, so a useful hint is to always use parentheses to force exactly the order you want when exponentiation is involved:

```
>>> 2 ** 3 ** 2      # The right-most ** operator gets done first!
512
>>> (2 ** 3) ** 2    # Use parentheses to force the order you want!
64
```

The immediate mode command prompt of Python is great for exploring and experimenting with expressions like this.

## 2.9 Operations on strings

In general, you cannot perform mathematical operations on strings, even if the strings look like numbers. The following are illegal (assuming that `message` has type string):

```
>>> message - 1      # Error
>>> "Hello" / 123     # Error
>>> message * "Hello" # Error
>>> "15" + 2          # Error
```

Interestingly, the `+` operator does work with strings, but for strings, the `+` operator represents **concatenation**, not addition. Concatenation means joining the two operands by linking them end-to-end. For example:

```
1 fruit = "banana"
2 baked_good = " nut bread"
3 print(fruit + baked_good)
```

The output of this program is `banana nut bread`. The space before the word `nut` is part of the string, and is necessary to produce the space between the concatenated strings.

The `*` operator also works on strings; it performs repetition. For example, `'Fun'*3` is `'FunFunFun'`. One of the operands has to be a string; the other has to be an integer.

On one hand, this interpretation of `+` and `*` makes sense by analogy with addition and multiplication. Just as  $4*3$  is equivalent to  $4+4+4$ , we expect `"Fun"*3` to be the same as `"Fun"+"Fun"+"Fun"`, and it is. On the other hand, there is a significant way in which string concatenation and repetition are different from integer addition and multiplication. Can you think of a property that addition and multiplication have that string concatenation and repetition do not?

## 2.10 Input

There is a built-in function in Python for getting input from the user:

```
1 name = input("Please enter your name: ")
```

The user of the program can enter the name and click *OK*, and when this happens the text that has been entered is returned from the `input` function, and in this case assigned to the variable `name`.

Even if you asked the user to enter their age, you would get back a string like `"17"`. It would be your job, as the programmer, to convert that string into a `int` or a `float`, using the `int` or `float` converter functions we saw earlier.

## 2.11 Composition

So far, we have looked at the elements of a program — variables, expressions, statements, and function calls — in isolation, without talking about how to combine them.

One of the most useful features of programming languages is their ability to take small building blocks and **compose** them into larger chunks.

For example, we know how to get the user to enter some input, we know how to convert the string we get into a float, we know how to write a complex expression, and we know how to print values. Let's put these together in a small four-step program that asks the user to input a value for the radius of a circle, and then computes the area of the circle from the formula

$$\text{Area} = \pi R^2$$

Firstly, we'll do the four steps one at a time:

```
1 response = input("What is your radius? ")
2 r = float(response)
3 area = 3.14159 * r**2
4 print("The area is ", area)
```

Now let's compose the first two lines into a single line of code, and compose the second two lines into another line of code.

```
1 r = float( input("What is your radius? ") )
2 print("The area is ", 3.14159 * r**2)
```

If we really wanted to be tricky, we could write it all in one statement:

```
1 print("The area is ", 3.14159*float(input("What is your radius?"))**2)
```

Such compact code may not be most understandable for humans, but it does illustrate how we can compose bigger chunks from our building blocks.

If you're ever in doubt about whether to compose code or fragment it into smaller steps, try to make it as simple as you can for the human to follow. My choice would be the first case above, with four separate steps.

## 2.12 The modulus operator

The **modulus operator** works on integers (and integer expressions) and gives the remainder when the first number is divided by the second. In Python, the modulus operator is a percent sign (%). The syntax is the same as for other operators. It has the same precedence as the multiplication operator.

```
>>> q = 7 // 3      # This is integer division operator
>>> print(q)
2
>>> r = 7 % 3
>>> print(r)
1
```

So 7 divided by 3 is 2 with a remainder of 1.

The modulus operator turns out to be surprisingly useful. For example, you can check whether one number is divisible by another—if `x % y` is zero, then `x` is divisible by `y`.

Also, you can extract the right-most digit or digits from a number. For example, `x % 10` yields the right-most digit of `x` (in base 10). Similarly `x % 100` yields the last two digits.

It is also extremely useful for doing conversions, say from seconds, to hours, minutes and seconds. So let's write a program to ask the user to enter some seconds, and we'll convert them into hours, minutes, and remaining seconds.

```
1 total_secs = int(input("How many seconds, in total?"))
2 hours = total_secs // 3600
3 secs_still_remaining = total_secs % 3600
4 minutes = secs_still_remaining // 60
5 secs_finally_remaining = secs_still_remaining % 60
6
7 print("Hrs=", hours, " mins=", minutes,
8       "secs=", secs_finally_remaining)
```

## 2.13 Glossary

**assignment statement** A statement that assigns a value to a name (variable). To the left of the assignment operator, `=`, is a name. To the right of the assignment token is an expression which is evaluated by the Python interpreter and then assigned to the name. The difference between the left and right hand sides of the assignment statement is often confusing to new programmers. In the following assignment:

```
number = number + 1
```

`number` plays a very different role on each side of the `=`. On the right it is a *value* and makes up part of the *expression* which will be evaluated by the Python interpreter before assigning it to the name on the left.

**assignment token** `=` is Python's assignment token. Do not confuse it with *equals*, which is an operator for comparing values.

**composition** The ability to combine simple expressions and statements into compound statements and expressions in order to represent complex computations concisely.

**concatenate** To join two strings end-to-end.

**data type** A set of values. The type of a value determines how it can be used in expressions. So far, the types you have seen are integers (`int`), floating-point numbers (`float`), and strings (`str`).

**evaluate** To simplify an expression by performing the operations in order to yield a single value.

**expression** A combination of variables, operators, and values that represents a single result value.

**float** A Python data type which stores *floating-point* numbers. Floating-point numbers are stored internally in two parts: a *base* and an *exponent*. When printed in the standard format, they look like decimal numbers. Beware of rounding errors when you use `floats`, and remember that they are only approximate values.

**floor division** An operator (denoted by the token `//`) that divides one number by another and yields an integer, or, if the result is not already an integer, it yields the next smallest integer.

**int** A Python data type that holds positive and negative whole numbers.

**keyword** A reserved word that is used by the compiler to parse program; you cannot use keywords like `if`, `def`, and `while` as variable names.

**modulus operator** An operator, denoted with a percent sign (`%`), that works on integers and yields the remainder when one number is divided by another.

**operand** One of the values on which an operator operates.

**operator** A special symbol that represents a simple computation like addition, multiplication, or string concatenation.

**rules of precedence** The set of rules governing the order in which expressions involving multiple operators and operands are evaluated.

**state snapshot** A graphical representation of a set of variables and the values to which they refer, taken at a particular instant during the program's execution.

**statement** An instruction that the Python interpreter can execute. So far we have only seen the assignment statement, but we will soon meet the `import` statement and the `for` statement.

**str** A Python data type that holds a string of characters.

**value** A number or string (or other things to be named later) that can be stored in a variable or computed in an expression.

**variable** A name that refers to a value.

**variable name** A name given to a variable. Variable names in Python consist of a sequence of letters (`a..z`, `A..Z`, and `_`) and digits (`0..9`) that begins with a letter. In best programming practice, variable names should be chosen so that they describe their use in the program, making the program *self documenting*.

## 2.14 Exercises

1. Take the sentence: *All work and no play makes Jack a dull boy*. Store each word in a separate variable, then print out the sentence on one line using `print`.
2. Add parenthesis to the expression `6 * 1 - 2` to change its value from 4 to -6.
3. Place a comment before a line of code that previously worked, and record what happens when you rerun the program.
4. Start the Python interpreter and enter `bruce + 4` at the prompt. This will give you an error:

```
NameError: name 'bruce' is not defined
```

Assign a value to `bruce` so that `bruce + 4` evaluates to 10.

5. The formula for computing the final amount if one is earning compound interest is given on Wikipedia as



$$A = P \left( 1 + \frac{r}{n} \right)^{nt}$$

Where,

- $P$  = principal amount (initial investment)
- $r$  = annual nominal interest rate (as a decimal)
- $n$  = number of times the interest is compounded per year
- $t$  = number of years

Here,  $P$  is the principal amount (the amount that the interest is provided on),  $n$  the frequency that the interest is paid out (per year), and  $r$  the interest rate. The number of years that the interest is calculated for is  $t$ . Write a program that replaces these letters with something a bit more human-readable, and calculate the interest for some varying amounts of money at realistic interest rates such as 1%, and -0.05%. When you have that working, ask the user for the value of some of these variables and do the calculation.

6. Evaluate the following numerical expressions in your head, then use the Python interpreter to check your results:

- (a) `>>> 5 % 2`
- (b) `>>> 9 % 5`
- (c) `>>> 15 % 12`
- (d) `>>> 12 % 15`
- (e) `>>> 6 % 6`
- (f) `>>> 0 % 7`
- (g) `>>> 7 % 0`

What happened with the last example? Why? If you were able to correctly anticipate the computer's response in all but the last one, it is time to move on. If not, take time now to make up examples of your own. Explore the modulus operator until you are confident you understand how it works.

- 7. You look at the clock and it is exactly 2pm. You set an alarm to go off in 51 hours. At what time does the alarm go off? (Hint: you could count on your fingers, but this is not what we're after. If you are tempted to count on your fingers, change the 51 to 5100.)
- 8. Write a Python program to solve the general version of the above problem. Ask the user for the time now (in hours), and ask for the number of hours to wait. Your program should output what the time will be on the clock when the alarm goes off.

### 3.1 Hello, little turtles!

There are many *modules* in Python that provide very powerful features that we can use in our own programs. Some of these can send email, or fetch web pages. The one we'll look at in this chapter allows us to create turtles and get them to draw shapes and patterns.

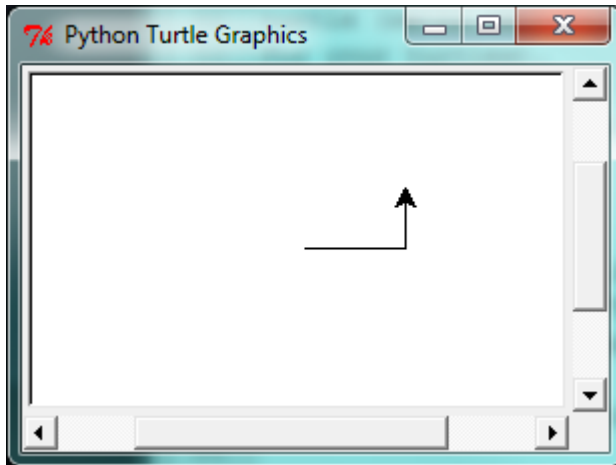
The turtles are fun, but the real purpose of the chapter is to teach ourselves a little more Python, and to develop our theme of *computational thinking*, or *thinking like a computer scientist*. Most of the Python covered here will be explored in more depth later.

#### 3.1.1 Our first turtle program

Let's write a couple of lines of Python program to create a new turtle and start drawing a rectangle. (We'll call the variable that refers to our first turtle `alex`, but we can choose another name if we follow the naming rules from the previous chapter).

```
1  import turtle                # Allows us to use turtles
2  window = turtle.Screen()     # Creates a playground for turtles
3  alex = turtle.Turtle()       # Create a turtle, assign to alex
4
5  alex.forward(50)              # Tell alex to move forward by 50 units
6  alex.left(90)                 # Tell alex to turn by 90 degrees
7  alex.forward(30)              # Complete the second side of a rectangle
8
9  window.mainloop()            # Wait for user to close window
```

When we run this program, a new window pops up:



Here are a couple of things we'll need to understand about this program.

The first line tells Python to load a module named `turtle`. That module brings us two new types that we can use: the `Turtle` type, and the `Screen` type. The dot notation `turtle.Turtle` means “*The Turtle type that is defined within the turtle module*”. (Remember that Python is case sensitive, so the module name, with a lowercase `t`, is different from the type `Turtle`.)

We then create and open what it calls a screen (we would prefer to call it a window), which we assign to variable `window`. Every window contains a **canvas**, which is the area inside the window on which we can draw.

In line 3 we create a turtle. The variable `alex` is made to refer to this turtle.

So these first three lines have set things up, we're ready to get our turtle to draw on our canvas.

In lines 5-7, we instruct the **object** `alex` to move, and to turn. We do this by **invoking**, or activating, `alex`'s **methods** — these are the instructions that all turtles know how to respond to.

The last line plays a part too: the `window` variable refers to the window shown above. When we invoke its `mainloop` method, it enters a state where it waits for events (like keypresses, or mouse movement and clicks). The program will terminate when the user closes the window.

An object can have various methods — things it can do — and it can also have **attributes** — (sometimes called *properties*). For example, each turtle has a *color* attribute. The method invocation `alex.color("red")` will make `alex` red, and drawing will be red too. (Note the word *color* is spelled the American way!)

The color of the turtle, the width of its pen, the position of the turtle within the window, which way it is facing, and so on are all part of its current **state**. Similarly, the window object has a background color, and some text in the title bar, and a size and position on the screen. These are all part of the state of the window object.

Quite a number of methods exist that allow us to modify the turtle and the window objects. We'll just show a couple. In this program we've only commented those lines that are different from the previous example (and we've used a different variable name for this turtle):

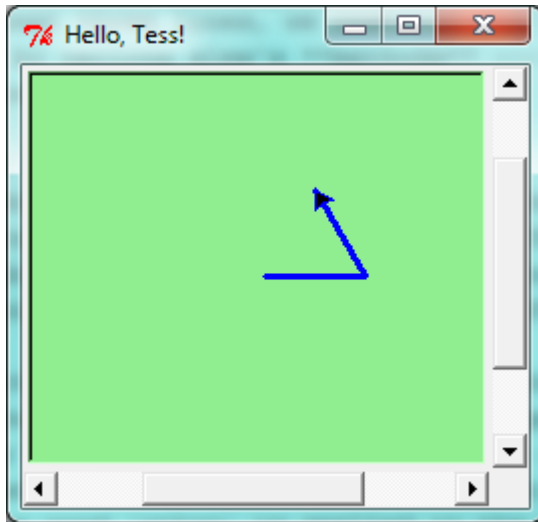
```
1  import turtle
2  window = turtle.Screen()
3  window.bgcolor("lightgreen")      # Set the window background color
4  window.title("Hello, Tess!")     # Set the window title
5
6  tess = turtle.Turtle()
7  tess.color("blue")               # Tell tess to change her color
8  tess.pensize(3)                 # Tell tess to set her pen width
9
10 tess.forward(50)
11 tess.left(120)
```

(continues on next page)

(continued from previous page)

```
12 tess.forward(50)
13
14 window.mainloop()
```

When we run this program, this new window pops up, and will remain on the screen until we close it.



---

#### Extend this program ...

1. Modify this program so that before it creates the window, it prompts the user to enter the desired background color. It should store the user's responses in a variable, and modify the color of the window according to the user's wishes. (Hint: you can find a list of permitted color names at <http://www.tcl.tk/man/tcl8.4/TkCmd/colors.htm>. It includes some quite unusual ones, like "peach puff" and "HotPink".)
  2. Do similar changes to allow the user, at runtime, to set `tess`' color.
- 

### 3.1.2 Instances — a herd of turtles

Just like we can have many different integers in a program, we can have many turtles. Each of them is called an **instance**. Each instance has its own attributes and methods — so `alex` might draw with a thin black pen and be at some position, while `tess` might be going in her own direction with a fat pink pen.

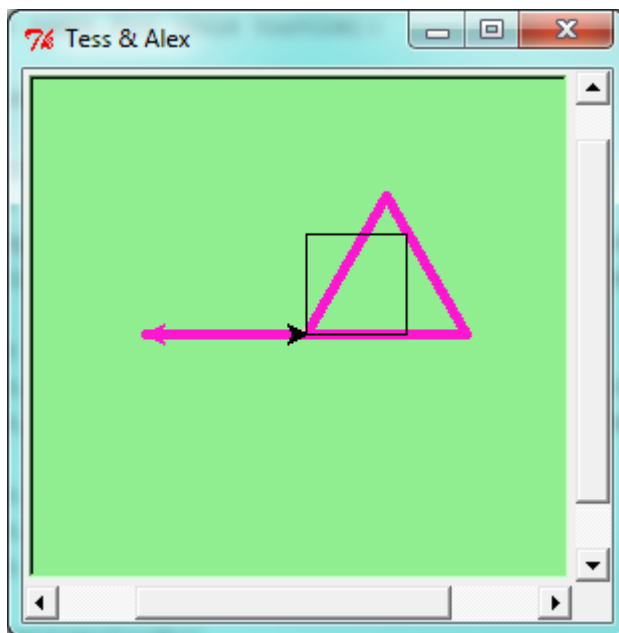
```
1  import turtle
2  window = turtle.Screen()           # Set up the window and its attributes
3  window.bgcolor("lightgreen")
4  window.title("Tess & Alex")
5
6  tess = turtle.Turtle()             # Create tess and set some attributes
7  tess.color("hotpink")
8  tess.pensize(5)
9
10 alex = turtle.Turtle()             # Create alex
11
12 tess.forward(80)                   # Make tess draw equilateral triangle
13 tess.left(120)
14 tess.forward(80)
```

(continues on next page)

(continued from previous page)

```
15 tess.left(120)
16 tess.forward(80)
17 tess.left(120)           # Complete the triangle
18
19 tess.right(180)          # Turn tess around
20 tess.forward(80)        # Move her away from the origin
21
22 alex.forward(50)         # Make alex draw a square
23 alex.left(90)
24 alex.forward(50)
25 alex.left(90)
26 alex.forward(50)
27 alex.left(90)
28 alex.forward(50)
29 alex.left(90)
30
31 window.mainloop()
```

Here is what happens when alex completes his rectangle, and tess completes her triangle:



Here are some *How to think like a computer scientist* observations:

- There are 360 degrees in a full circle. If we add up all the turns that a turtle makes, *no matter what steps occurred between the turns*, we can easily figure out if they add up to some multiple of 360. This should convince us that alex is facing in exactly the same direction as he was when he was first created. (Geometry conventions have 0 degrees facing East, and that is the case here too!)
- We could have left out the last turn for alex, but that would not have been as satisfying. If we're asked to draw a closed shape like a square or a rectangle, it is a good idea to complete all the turns and to leave the turtle back where it started, facing the same direction as it started in. This makes reasoning about the program and composing chunks of code into bigger programs easier for us humans!
- We did the same with tess: she drew her triangle, and turned through a full 360 degrees. Then we turned her around and moved her aside. Even the blank line 18 is a hint about how the programmer's *mental chunking* is working: in big terms, tess' movements were chunked as "draw the triangle" (lines 12-17) and then "move away from the origin" (lines 19 and 20).

- One of the key uses for comments is to record our mental chunking, and big ideas. They're not always explicit in the code.
- And, uh-huh, two turtles may not be enough for a herd. But the important idea is that the turtle module gives us a kind of factory that lets us create as many turtles as we need. Each instance has its own state and behaviour.

### 3.1.3 The for loop

When we drew the square, it was quite tedious. We had to explicitly repeat the steps of moving and turning four times. If we were drawing a hexagon, or an octagon, or a polygon with 42 sides, it would have been worse.

So a basic building block of all programs is to be able to repeat some code, over and over again.

Python's **for** loop solves this for us. Let's say we have some friends, and we'd like to send them each an email inviting them to our party. We don't quite know how to send email yet, so for the moment we'll just print a message for each friend:

```
1  for friend in ["Joe", "Zoe", "Zuki", "Thandi", "Paris"]:  
2      invite = "Hi " + friend + ". Please come to my party!"  
3      print(invite)  
4  # more code can follow here ...
```

When we run this, the output looks like this:

```
Hi Joe. Please come to my party!  
Hi Zoe. Please come to my party!  
Hi Zuki. Please come to my party!  
Hi Thandi. Please come to my party!  
Hi Paris. Please come to my party!
```

- The variable `friend` in the `for` statement at line 1 is called the **loop variable**. We could have chosen any other variable name instead, such as `broccoli`: the computer doesn't care.
- Lines 2 and 3 are the **loop body**. The loop body is always indented. The indentation determines exactly what statements are "in the body of the loop".
- On each *iteration* or *pass* of the loop, first a check is done to see if there are still more items to be processed. If there are none left (this is called the **terminating condition** of the loop), the loop has finished. Program execution continues at the next statement after the loop body, (e.g. in this case the next statement below the comment in line 4).
- If there are items still to be processed, the loop variable is updated to refer to the next item in the list. This means, in this case, that the loop body is executed here 7 times, and each time `friend` will refer to a different friend.
- At the end of each execution of the body of the loop, Python returns to the `for` statement, to see if there are more items to be handled, and to assign the next one to `friend`.

### 3.1.4 Flow of Execution of the for loop

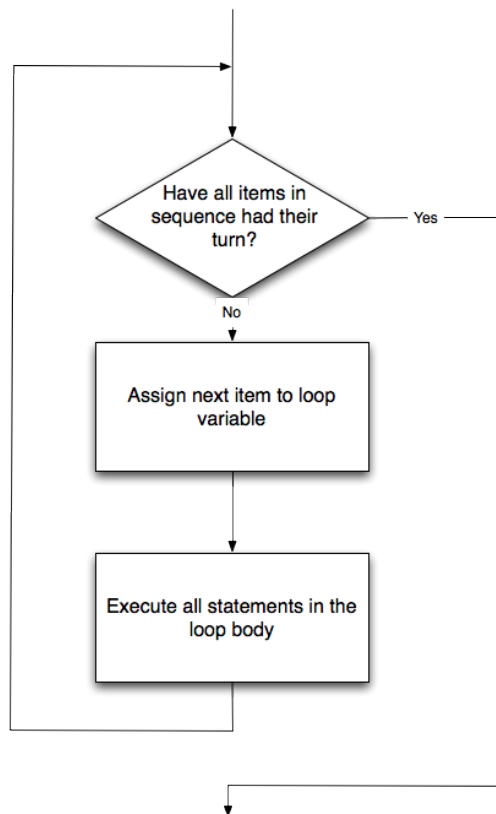
As a program executes, the interpreter always keeps track of which statement is about to be executed. We call this the **control flow**, or the **flow of execution** of the program. When humans execute programs, they often use their finger to point to each statement in turn. So we could think of control flow as "Python's moving finger".

Control flow until now has been strictly top to bottom, one statement at a time. The `for` loop changes this.

---

#### Flowchart of a for loop

Control flow is often easy to visualize and understand if we draw a flowchart. This shows the exact steps and logic of how the `for` statement executes.



### 3.1.5 The loop simplifies our turtle program

To draw a square we'd like to do the same thing four times — move the turtle, and turn. We previously used 8 lines to have `alex` draw the four sides of a square. This does exactly the same, but using just three lines:

```
1 for i in [0,1,2,3]:
2     alex.forward(50)
3     alex.left(90)
```

Some observations:

- While “saving some lines of code” might be convenient, it is not the big deal here. What is much more important is that we’ve found a “repeating pattern” of statements, and reorganized our program to repeat the pattern. Finding the chunks and somehow getting our programs arranged around those chunks is a vital skill in computational thinking.
- The values `[0,1,2,3]` were provided to make the loop body execute 4 times. We could have used any four values, but these are the conventional ones to use. In fact, they are so popular that Python gives us special built-in `range` objects:

```
1 for i in range(4):
2     # Executes the body with i = 0, then 1, then 2, then 3
3 for x in range(10):
4     # Sets x to each of ... [0, 1, 2, 3, 4, 5, 6, 7, 8, 9]
```

- Since we do not need or use the variable `i` in this case, we could replace it with `_`, although this is not important for the program flow, it is good style.
- Computer scientists like to count from 0!
- `range` can deliver a sequence of values to the loop variable in the `for` loop. They start at 0, and in these cases do not include the 4 or the 10.
- Our little trick earlier to make sure that `alex` did the final turn to complete 360 degrees has paid off: if we had not done that, then we would not have been able to use a loop for the fourth side of the square. It would have become a “special case”, different from the other sides. When possible, we’d much prefer to make our code fit a general pattern, rather than have to create a special case.

So to repeat something four times, a good Python programmer would do this:

```
1 for _ in range(4):
2     alex.forward(50)
3     alex.left(90)
```

By now you should be able to see how to change our previous program so that `tess` can also use a `for` loop to draw her equilateral triangle.

But now, what would happen if we made this change?

```
1 for color in ["yellow", "red", "purple", "blue"]:
2     alex.color(color)
3     alex.forward(50)
4     alex.left(90)
```

A variable can also be assigned a value that is a list. So lists can also be used in more general situations, not only in the `for` loop. The code above could be rewritten like this:

```
1 # Assign a list to a variable
2 colors = ["yellow", "red", "purple", "blue"]
3 for color in colors:
4     alex.color(color)
5     alex.forward(50)
6     alex.left(90)
```

- Notice the difference between the method `alex.color`, which is “part of” the instance `alex`, and the variable `color`, which is “part of” the main body of your program.

### 3.1.6 A few more turtle methods and tricks

Turtle methods can use negative angles or distances. So `tess.forward(-100)` will move `tess` backwards, and `tess.left(-30)` turns her to the right. Additionally, because there are 360 degrees in a circle, turning 30 to the left will get `tess` facing in the same direction as turning 330 to the right! (The on-screen animation will differ, though — you will be able to tell if `tess` is turning clockwise or counter-clockwise!)

This suggests that we don’t need both a left and a right turn method — we could be minimalists, and just have one method. There is also a *backward* method. (If you are very nerdy, you might enjoy saying `alex.backward(-100)` to move `alex` forward!)

Part of *thinking like a scientist* is to understand more of the structure and rich relationships in our field. So revising a few basic facts about geometry and number lines, and spotting the relationships between left, right, backward, forward, negative and positive distances or angles values is a good start if we’re going to play with turtles.

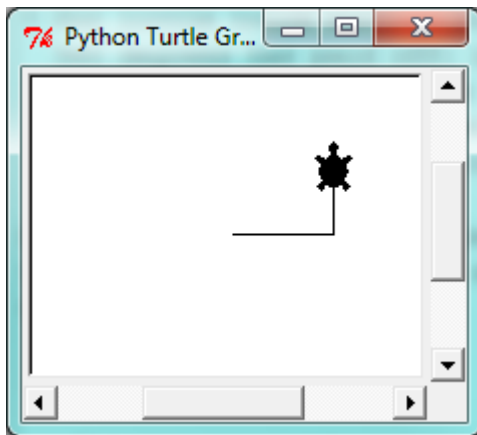
A turtle’s pen can be picked up or put down. This allows us to move a turtle to a different place without drawing a line. The methods are



```
1 alex.penup()
2 alex.forward(100)      # This moves alex, but no line is drawn
3 alex.pendown()
```

Every turtle can have its own shape. The ones available “out of the box” are arrow, blank, circle, classic, square, triangle, turtle.

```
1 alex.shape("turtle")
```



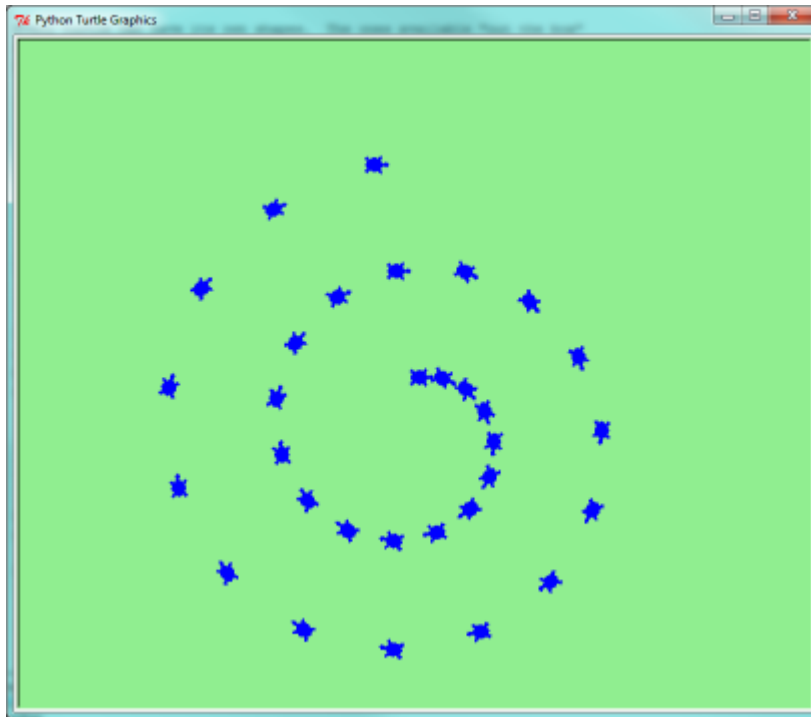
We can speed up or slow down the turtle’s animation speed. (Animation controls how quickly the turtle turns and moves forward). Speed settings can be set between 1 (slowest) to 10 (fastest). But if we set the speed to 0, it has a special meaning — turn off animation and go as fast as possible.

```
1 alex.speed(10)
```

A turtle can “stamp” its footprint onto the canvas, and this will remain after the turtle has moved somewhere else. Stamping works, even when the pen is up.

Let’s do an example that shows off some of these new features:

```
1 import turtle
2 window = turtle.Screen()
3 window.bgcolor("lightgreen")
4 tess = turtle.Turtle()
5 tess.shape("turtle")
6 tess.color("blue")
7
8 tess.penup()                # This is new
9 size = 20
10 for _ in range(30):
11     tess.stamp()           # Leave an impression on the canvas
12     size = size + 3        # Increase the size on every iteration
13     tess.forward(size)     # Move tess along
14     tess.right(24)         # ... and turn her
15
16 window.mainloop()
```



Be careful now! How many times was the body of the loop executed? How many turtle images do we see on the screen? All except one of the shapes we see on the screen here are footprints created by `stamp`. But the program still only has *one* turtle instance — can you figure out which one here is the real `tess`? (Hint: if you're not sure, write a new line of code after the `for` loop to change `tess`' color, or to put her pen down and draw a line, or to change her shape, etc.)

## 3.2 Conditionals

Programs get really interesting when we can test conditions and change the program behaviour depending on the outcome of the tests. That's what this part is about.

### 3.2.1 Boolean values and expressions

A *Boolean* value is either true or false. It is named after the British mathematician, George Boole, who first formulated *Boolean algebra* — some rules for reasoning about and combining these values. This is the basis of all modern computer logic.

In Python, the two Boolean values are `True` and `False` (the capitalization must be exactly as shown), and the Python type is `bool`.

```
>>> type(True)
<class 'bool'>
>>> type(true)
Traceback (most recent call last):
  File "<interactive input>", line 1, in <module>
NameError: name 'true' is not defined
```

A **Boolean expression** is an expression that evaluates to produce a result which is a Boolean value. For example, the operator `==` tests if two values are equal. It produces (or *yields*) a Boolean value:

```
>>> 5 == (3 + 2)    # Is five equal 5 to the result of 3 + 2?
True
>>> 5 == 6
False
>>> j = "hel"
>>> j + "lo" == "hello"
True
```

In the first statement, the two operands evaluate to equal values, so the expression evaluates to `True`; in the second statement, 5 is not equal to 6, so we get `False`.

The `==` operator is one of six common **comparison operators** which all produce a `bool` result; here are all six:

```
x == y          # Produce True if ... x is equal to y
x != y          # ... x is not equal to y
x > y           # ... x is greater than y
x < y           # ... x is less than y
x >= y          # ... x is greater than or equal to y
x <= y          # ... x is less than or equal to y
```

Although these operations are probably familiar, the Python symbols are different from the mathematical symbols. A common error is to use a single equal sign (`=`) instead of a double equal sign (`==`). Remember that `=` is an assignment operator and `==` is a comparison operator. Also, there is no such thing as `=<` or `=>`.

Like any other types we've seen so far, Boolean values can be assigned to variables, printed, etc.

```
>>> age = 19
>>> old_enough_to_get_driving_licence = age >= 18
>>> print(old_enough_to_get_driving_licence)
True
>>> type(old_enough_to_get_driving_licence)
<class 'bool'>
```

### 3.2.2 Logical operators

There are three **logical operators**, `and`, `or`, and `not`, that allow us to build more complex Boolean expressions from simpler Boolean expressions. The semantics (meaning) of these operators is similar to their meaning in English. For example, `x > 0` and `x < 10` produces `True` only if `x` is greater than 0 *and* at the same time, `x` is less than 10.

`n % 2 == 0` or `n % 3 == 0` is `True` if *either* of the conditions is `True`, that is, if the number `n` is divisible by 2 *or* it is divisible by 3. (What do you think happens if `n` is divisible by both 2 and by 3 at the same time? Will the expression yield `True` or `False`? Try it in your Python interpreter.)

Finally, the `not` operator negates a Boolean value, so `not (x > y)` is `True` if `(x > y)` is `False`, that is, if `x` is less than or equal to `y`. In other words: `not True` is `False`, and `not False` is `True`.

The expression on the left of the `or` operator is evaluated first: if the result is `True`, Python does not (and need not) evaluate the expression on the right — this is called *short-circuit evaluation*. Similarly, for the `and` operator, if the expression on the left yields `False`, Python does not evaluate the expression on the right.

So there are no unnecessary evaluations.

### 3.2.3 Truth Tables

A truth table is a small table that allows us to list all the possible inputs, and to give the results for the logical operators. Because the `and` and `or` operators each have two operands, there are only four rows in a truth table that describes the

semantics of `and`.

a	b	a and b
False	False	False
False	True	False
True	False	False
True	True	True

In a Truth Table, we sometimes use T and F as shorthand for the two Boolean values: here is the truth table describing `or`:

a	b	a or b
F	F	F
F	T	T
T	F	T
T	T	T

The third logical operator, `not`, only takes a single operand, so its truth table only has two rows:

a	not a
F	T
T	F

### 3.2.4 Simplifying Boolean Expressions

A set of rules for simplifying and rearranging expressions is called an *algebra*. For example, we are all familiar with school algebra rules, such as:

```
n * 0 == 0
```

Here we see a different algebra — the *Boolean* algebra — which provides rules for working with Boolean values.

First, the `and` operator:

```
x and False == False
False and x == False
y and x == x and y
x and True == x
True and x == x
x and x == x
```

Here are some corresponding rules for the `or` operator:

```
x or False == x
False or x == x
y or x == x or y
x or True == True
True or x == True
x or x == x
```

Two `not` operators cancel each other:

```
not (not x) == x
```

### 3.2.5 Conditional execution

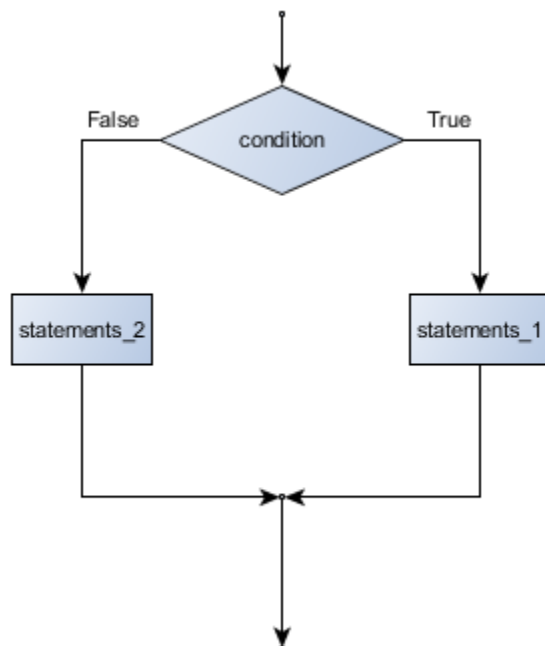
In order to write useful programs, we almost always need the ability to check conditions and change the behavior of the program accordingly. **Conditional statements** give us this ability. The simplest form is the **if** statement:

```
1 if x % 2 == 0:
2     print(x, " is even.")
3     print("Did you know that 2 is the only even number that is prime?")
4 else:
5     print(x, " is odd.")
6     print("Did you know that multiplying two odd numbers " +
7           "always gives an odd result?")
```

The Boolean expression after the **if** statement is called the **condition**. If it is true, then all the indented statements get executed. If not, then all the statements indented under the **else** clause get executed.

---

#### Flowchart of an if statement with an else clause



---

The syntax for an **if** statement looks like this:

```
1 if <BOOLEAN EXPRESSION>:
2     <STATEMENTS_1>           # Executed if condition evaluates to True
3 else:
4     <STATEMENTS_2>           # Executed if condition evaluates to False
```

As with the function definition from the next chapter and other compound statements like **for**, the **if** statement consists of a header line and a body. The header line begins with the keyword **if** followed by a *Boolean expression* and ends with a colon (:).

The indented statements that follow are called a **block**. The first unindented statement marks the end of the block.

Each of the statements inside the first block of statements are executed in order if the Boolean expression evaluates to `True`. The entire first block of statements is skipped if the Boolean expression evaluates to `False`, and instead all the statements indented under the `else` clause are executed.

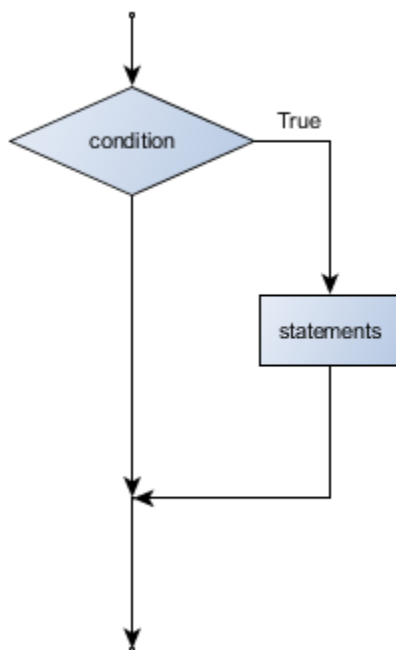
There is no limit on the number of statements that can appear under the two clauses of an `if` statement, but there has to be at least one statement in each block. Occasionally, it is useful to have a section with no statements (usually as a place keeper, or scaffolding, for code we haven't written yet). In that case, we can use the `pass` statement, which does nothing except act as a placeholder.

```
1  if True:           # This is always True,  
2      pass          # so this is always executed, but it does nothing  
3  else:  
4      pass          # And this is never executed
```

### 3.2.6 Omitting the `else` clause

---

Flowchart of an `if` statement with no `else` clause



---

Another form of the `if` statement is one in which the `else` clause is omitted entirely. In this case, when the condition evaluates to `True`, the statements are executed, otherwise the flow of execution continues to the statement after the `if`.

```
1  if x < 0:  
2      print("The negative number ", x, " is not valid here.")  
3      x = 42  
4      print("I've decided to use the number 42 instead.")  
5  
6  print("The square root of ", x, "is", math.sqrt(x))
```

In this case, the `print` function that outputs the square root is the one after the `if` — not because we left a blank line, but because of the way the code is indented. Note too that the function call `math.sqrt(x)` will give an error unless

we have an `import math` statement, usually placed near the top of our script.

---

### Python terminology

Python documentation sometimes uses the term **suite** of statements to mean what we have called a *block* here. They mean the same thing, and since most other languages and computer scientists use the word *block*, we'll stick with that.

Notice too that `else` is not a statement. The `if` statement has two *clauses*, one of which is the (optional) `else` clause.

---

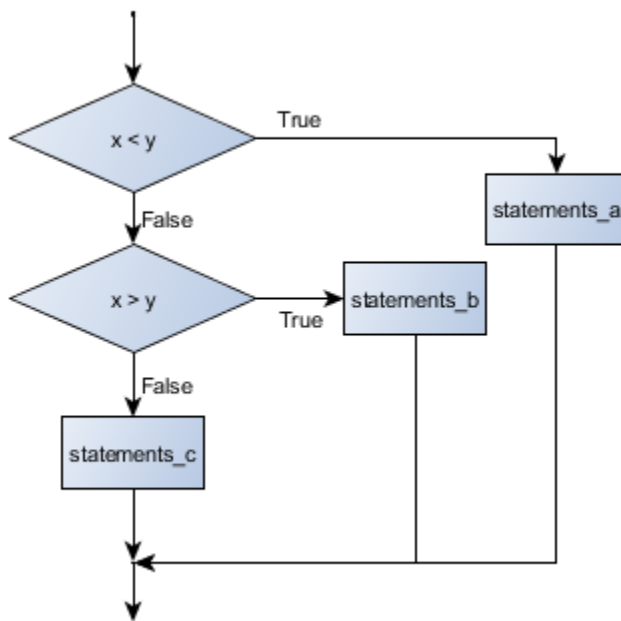
### 3.2.7 Chained conditionals

Sometimes there are more than two possibilities and we need more than two branches. One way to express a computation like that is a **chained conditional**:

```
1  if x < y:
2      <STATEMENTS_A>
3  elif x > y:
4      <STATEMENTS_B>
5  else:                # x == y
6      <STATEMENTS_C>
```

---

#### Flowchart of this chained conditional



---

`elif` is an abbreviation of `else if`. Again, exactly one branch will be executed. There is no limit of the number of `elif` statements but only a single (and optional) final `else` statement is allowed and it must be the last branch in the statement:

```
1 if choice == "a":
2     function_one()
3 elif choice == "b":
4     function_two()
5 elif choice == "c":
6     function_three()
7 else:
8     print("Invalid choice.")
```

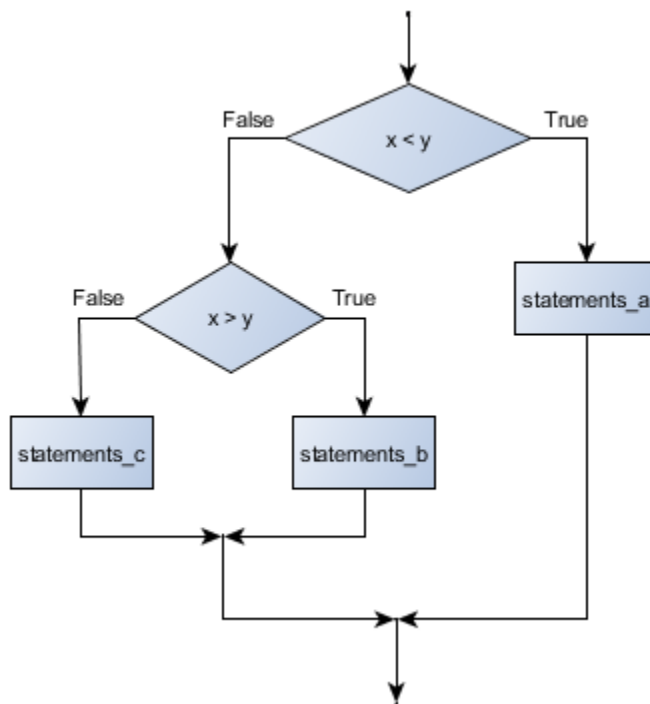
Each condition is checked in order. If the first is false, the next is checked, and so on. If one of them is true, the corresponding branch executes, and the statement ends. Even if more than one condition is true, only the first true branch executes.

### 3.2.8 Nested conditionals

One conditional can also be **nested** within another. (It is the same theme of composability, again!) We could have written the previous example as follows:

---

Flowchart of this nested conditional



```
1 if x < y:
2     <STATEMENTS_A>
3 else:
4     if x > y:
5         <STATEMENTS_B>
6     else:
7         <STATEMENTS_C>
```



The outer conditional contains two branches. The second branch contains another `if` statement, which has two branches of its own. Those two branches could contain conditional statements as well.

Although the indentation of the statements makes the structure apparent, nested conditionals very quickly become very difficult to read. In general, it is a good idea to avoid them when we can.

Logical operators often provide a way to simplify nested conditional statements. For example, we can rewrite the following code using a single conditional:

```
1 if 0 < x:                # Assume x is an int here
2     if x < 10:
3         print("x is a positive single digit.")
```

The `print` function is called only if we make it past both the conditionals, so instead of the above which uses two `if` statements each with a simple condition, we could make a more complex condition using the `and` operator. Now we only need a single `if` statement:

```
1 if 0 < x and x < 10:
2     print("x is a positive single digit.")
```

In this case there is a third option:

```
1 if 0 < x < 10:
2     print("x is a positive single digit.")
```

### 3.2.9 Logical opposites

Each of the six relational operators has a logical opposite: for example, suppose we can get a driving licence when our age is greater or equal to 18, we can *not* get the driving licence when we are less than 18.

Notice that the opposite of `>=` is `<`.

operator	logical opposite
<code>==</code>	<code>!=</code>
<code>!=</code>	<code>==</code>
<code>&lt;</code>	<code>&gt;=</code>
<code>&lt;=</code>	<code>&gt;</code>
<code>&gt;</code>	<code>&lt;=</code>
<code>&gt;=</code>	<code>&lt;</code>

Understanding these logical opposites allows us to sometimes get rid of `not` operators. `not` operators are often quite difficult to read in computer code, and our intentions will usually be clearer if we can eliminate them.

For example, if we wrote this Python:

```
1 if not (age >= 18):
2     print("Hey, you're too young to get a driving licence!")
```

it would probably be clearer to use the simplification laws, and to write instead:

```
1 if age < 18:
2     print("Hey, you're too young to get a driving licence!")
```

Two powerful simplification laws (called de Morgan's laws) that are often helpful when dealing with complicated Boolean expressions are:

```
(not (x and y)) == ((not x) or (not y))
(not (x or y))  == ((not x) and (not y))
```

For example, suppose we can slay the dragon only if our magic lightsabre sword is charged to 90% or higher, and we have 100 or more energy units in our protective shield. We find this fragment of Python code in the game:

```
1 if not (sword_charge >= 0.90 and shield_energy >= 100):
2     print("Your attack has no effect, the dragon fries you to a crisp!")
3 else:
4     print("The dragon crumples in a heap. You rescue the gorgeous princess!")
```

de Morgan's laws together with the logical opposites would let us rework the condition in a (perhaps) easier to understand way like this:

```
1 if sword_charge < 0.90 or shield_energy < 100:
2     print("Your attack has no effect, the dragon fries you to a crisp!")
3 else:
4     print("The dragon crumples in a heap. You rescue the gorgeous princess!")
```

We could also get rid of the not by swapping around the then and else parts of the conditional. So here is a third version, also equivalent:

```
1 if sword_charge >= 0.90 and shield_energy >= 100:
2     print("The dragon crumples in a heap. You rescue the gorgeous princess!")
3 else:
4     print("Your attack has no effect, the dragon fries you to a crisp!")
```

To improve readability, there is this fourth version:

```
1 sword_check = sword_charge >= 0.90
2 shield_check = shield_energy >= 100
3
4 if sword_check and shield_check:
5     print("The dragon crumples in a heap. You rescue the gorgeous princess!")
6 else:
7     print("Your attack has no effect, the dragon fries you to a crisp!")
```

This version is probably the best of the four, because it very closely matches the initial English statement. Clarity of our code (for other humans), and making it easy to see that the code does what was expected should always be highest priority.

As our programming skills develop we'll find we have more than one way to solve any problem. So good programs are *designed*. We make choices that favour clarity, simplicity, and elegance. The job title *software architect* says a lot about what we do — we are *architects* who engineer our products to balance beauty, functionality, simplicity and clarity in our creations.

---

**Tip:** Once our program works, we should play around a bit trying to polish it up. Write good comments. Think about whether the code would be clearer with different variable names. Could we have done it more elegantly? Should we rather use a function? Can we simplify the conditionals?

We think of our code as our creation, our work of art! We make it great.

---

## 3.3 Iteration

Computers are often used to automate repetitive tasks. Repeating identical or similar tasks without making errors is something that computers do well and people do poorly.

Repeated execution of a set of statements is called **iteration**. Because iteration is so common, Python provides several language features to make it easier. We've already seen the `for` statement. This is the form of iteration you'll likely be using most often. But here we're going to look at the `while` statement — another way to have your program do iteration, useful in slightly different circumstances.

Before we do that, let's just review a few ideas...

### 3.3.1 Assignment

As we have mentioned previously, it is legal to make more than one assignment to the same variable. A new assignment makes an existing variable refer to a new value (and stop referring to the old value).

```
1  airtime_remaining = 15
2  print(airtime_remaining)
3  airtime_remaining = 7
4  print(airtime_remaining)
```

The output of this program is:

```
15
7
```

because the first time `airtime_remaining` is printed, its value is 15, and the second time, its value is 7.

It is especially important to distinguish between an assignment statement and a Boolean expression that tests for equality. Because Python uses the equal token (`=`) for assignment, it is tempting to interpret a statement like `a = b` as a Boolean test. Unlike mathematics, it is not! Remember that the Python token for the equality operator is `==`.

Note too that an equality test is symmetric, but assignment is not. For example, if `a == 7` then `7 == a`. But in Python, the statement `a = 7` is legal and `7 = a` is not.

In Python, an assignment statement can make two variables equal, but because further assignments can change either of them, they don't have to stay that way:

```
1  a = 5
2  b = a      # After executing this line, a and b are now equal
3  a = 3      # After executing this line, a and b are no longer equal
```

The third line changes the value of `a` but does not change the value of `b`, so they are no longer equal. (In some programming languages, a different symbol is used for assignment, such as `<-` or `:=`, to avoid confusion. Some people also think that *variable* was an unfortunate word to choose, and instead we should have called them *assignables*. Python chooses to follow common terminology and token usage, also found in languages like C, C++, Java, and C#, so we use the tokens `=` for assignment, `==` for equality, and we talk of *variables*.

### 3.3.2 Updating variables

When an assignment statement is executed, the right-hand side expression (i.e. the expression that comes after the assignment token) is evaluated first. This produces a *value*. Then the assignment is made, so that the variable on the left-hand side now refers to the new value.

One of the most common forms of assignment is an update, where the new value of the variable depends on its old value. Deduct 40 cents from my airline balance, or add one run to the scoreboard.

```
1 n = 5
2 n = 3 * n + 1
```

Line 2 means *get the current value of  $n$ , multiply it by three and add one, and assign the answer to  $n$ , thus making  $n$  refer to the value*. So after executing the two lines above,  $n$  will point/refer to the integer 16.

If you try to get the value of a variable that has never been assigned to, you'll get an error:

```
>>> w = x + 1
Traceback (most recent call last):
  File "<interactive input>", line 1, in
NameError: name 'x' is not defined
```

Before you can update a variable, you have to **initialize** it to some starting value, usually with a simple assignment:

```
1 runs_scored = 0
2 ...
3 runs_scored = runs_scored + 1
```

Line 3 — updating a variable by adding 1 to it — is very common. It is called an **increment** of the variable; subtracting 1 is called a **decrement**. Sometimes programmers also talk about *bumping* a variable, which means the same as incrementing it by 1. This is commonly done with the `+=` operator.

```
1 runs_scored = 0
2 ...
3 runs_scored += 1
```

### 3.3.3 The `for` loop revisited

Recall that the `for` loop processes each item in a list. Each item in turn is (re-)assigned to the loop variable, and the body of the loop is executed. We saw this example before:

```
1 for friend in ["Joe", "Zoe", "Zuki", "Thandi", "Paris"]:
2     invite = "Hi " + friend + ". Please come to my party!"
3     print(invite)
```

Running through all the items in a list is called **traversing** the list, or **traversal**.

Let us write some code now to sum up all the elements in a list of numbers. Do this by hand first, and try to isolate exactly what steps you take. You'll find you need to keep some “running total” of the sum so far, either on a piece of paper, in your head, or in your calculator. Remembering things from one step to the next is precisely why we have variables in a program: so we'll need some variable to remember the “running total”. It should be initialized with a value of zero, and then we need to traverse the items in the list. For each item, we'll want to update the running total by adding the next number to it.

```
1 numbers = [5, 6, 32, 21, 9]
2 running_total = 0
3 for number in numbers:
4     running_total = running_total + number
5 print(running_total)
```

### 3.3.4 The `while` statement

Here is a fragment of code that demonstrates the use of the `while` statement:

```
1 while <CONDITION>:
2     <STATEMENT>

1 n = 6
2
3 current_sum = 0
4 i = 0
5 while i <= n:
6     current_sum += i
7     i += 1
8 print(current_sum)
```

You can almost read the `while` statement as if it were English. It means, while `i` is less than or equal to `n`, continue executing the body of the loop. Within the body, each time, increment `i`. When `i` passes `n`, return your accumulated sum. In other words: while `<CONDITION>` is `True`, `<STATEMENT>` is executed. Of course, this example could be written more concisely as `sum(range(n + 1))` because the function `sum` already exists.

More formally, here is precise flow of execution for a `while` statement:

- Evaluate the condition at line 5, yielding a value which is either `False` or `True`.
- If the value is `False`, exit the `while` statement and continue execution at the next statement (line 8 in this case).
- If the value is `True`, execute each of the statements in the body (lines 6 and 7) and then go back to the `while` statement at line 5.

The body consists of all of the statements indented below the `while` keyword.

Notice that if the loop condition is `False` the first time we get loop, the statements in the body of the loop are never executed.

The body of the loop should change the value of one or more variables so that eventually the condition becomes false and the loop terminates. Otherwise the loop will repeat forever, which is called an **infinite loop**.

In the case here, we can prove that the loop terminates because we know that the value of `n` is finite, and we can see that the value of `i` increments each time through the loop, so eventually it will have to exceed `n`. In other cases it is not so easy, maybe even impossible, to tell if the loop will ever terminate.

What you will notice here is that the `while` loop is more work for you — the programmer — than the equivalent `for` loop. When using a `while` loop one has to manage the loop variable yourself: give it an initial value, test for completion, and then make sure you change something in the body so that the loop terminates. By comparison, here is an equivalent snippet that uses `for` instead:

```
1 n = 6
2
3 current_sum = 0
4 for i in range(n+1):
5     current_sum += i
6 print(current_sum)
```

Notice the slightly tricky call to the `range` function — we had to add one onto `n`, because `range` generates its list up to but excluding the value you give it. It would be easy to make a programming mistake and overlook this.

So why have two kinds of loop if `for` looks easier? This next example shows a case where we need the extra power that we get from the `while` loop.

### 3.3.5 The Collatz $3n + 1$ sequence

Let's look at a simple sequence that has fascinated and foxed mathematicians for many years. They still cannot answer even quite simple questions about this.

The “computational rule” for creating the sequence is to start from some given  $n$ , and to generate the next term of the sequence from  $n$ , either by halving  $n$ , (whenever  $n$  is even), or else by multiplying it by three and adding 1. The sequence terminates when  $n$  reaches 1.

This Python snippet captures that algorithm:

```
1 n = 1027371
2
3 while n != 1:
4     print(n, end=", ")
5     if n % 2 == 0:           # n is even
6         n = n // 2
7     else:                   # n is odd
8         n = n * 3 + 1
9 print(n, end=".\n")
```

Notice first that the print function on line 4 has an extra argument `end=", "`. This tells the `print` function to follow the printed string with whatever the programmer chooses (in this case, a comma followed by a space), instead of ending the line. So each time something is printed in the loop, it is printed on the same output line, with the numbers separated by commas. The call to `print(n, end=".\n")` at line 9 after the loop terminates will then print the final value of  $n$  followed by a period and a newline character. (You'll cover the `\n` (newline character) later).

The condition for continuing with this loop is `n != 1`, so the loop will continue running until it reaches its termination condition, (i.e. `n == 1`).

Each time through the loop, the program outputs the value of  $n$  and then checks whether it is even or odd. If it is even, the value of  $n$  is divided by 2 using integer division. If it is odd, the value is replaced by `n * 3 + 1`.

Since  $n$  sometimes increases and sometimes decreases, there is no obvious proof that  $n$  will ever reach 1, or that the program terminates. For some particular values of  $n$ , we can prove termination. For example, if the starting value is a power of two, then the value of  $n$  will be even each time through the loop until it reaches 1. The previous example ends with such a sequence, starting with 16.

See if you can find a small starting number that needs more than a hundred steps before it terminates.

Particular values aside, the interesting question was first posed by a German mathematician called Lothar Collatz: the *Collatz conjecture* (also known as the  *$3n + 1$  conjecture*), is that this sequence terminates for *all* positive values of  $n$ . So far, no one has been able to prove it *or* disprove it! (A conjecture is a statement that might be true, but nobody knows for sure.)

Think carefully about what would be needed for a proof or disproof of the conjecture “*All positive integers will eventually converge to 1 using the Collatz rules*”. With fast computers we have been able to test every integer up to very large values, and so far, they have all eventually ended up at 1. But who knows? Perhaps there is some as-yet untested number which does not reduce to 1.

You'll notice that if you don't stop when you reach 1, the sequence gets into its own cyclic loop: 1, 4, 2, 1, 4, 2, 1, 4 ... So one possibility is that there might be other cycles that we just haven't found yet.

Wikipedia has an informative article about the Collatz conjecture. The sequence also goes under other names (Hailstone sequence, Wonderous numbers, etc.), and you'll find out just how many integers have already been tested by computer, and found to converge!

---

#### Choosing between `for` and `while`

Use a `for` loop if you know, before you start looping, the maximum number of times that you'll need to execute the body. For example, if you're traversing a list of elements, you know that the maximum number of loop iterations you can possibly need is "all the elements in the list". Or if you need to print the 12 times table, we know right away how many times the loop will need to run.

So any problem like "iterate this weather model for 1000 cycles", or "search this list of words", "find all prime numbers up to 10000" suggest that a `for` loop is best.

By contrast, if you are required to repeat some computation until some condition is met, and you cannot calculate in advance when (of if) this will happen, as we did in this  $3n + 1$  problem, you'll need a `while` loop.

We call the first case **definite iteration** — we know ahead of time some definite bounds for what is needed. The latter case is called **indefinite iteration** — we're not sure how many iterations we'll need — we cannot even establish an upper bound!

---

### 3.3.6 Tracing a program

To write effective computer programs, and to build a good conceptual model of program execution, a programmer needs to develop the ability to **trace** the execution of a computer program. Tracing involves becoming the computer and following the flow of execution through a sample program run, recording the state of all variables and any output the program generates after each instruction is executed.

To understand this process, let's trace the call to the collatz code above with  $n = 3$  from the previous section. At the start of the trace, we have a variable, `n`, with an initial value of 3. Since 3 is not equal to 1, the `while` loop body is executed. 3 is printed and `3 % 2 == 0` is evaluated. Since it evaluates to `False`, the `else` branch is executed and `3 * 3 + 1` is evaluated and assigned to `n`.

To keep track of all this as you hand trace a program, make a column heading on a piece of paper for each variable created as the program runs and another one for output. Our trace so far would look something like this:

n	output printed so far
--	-----
3	3,
10	

Since `10 != 1` evaluates to `True`, the loop body is again executed, and 10 is printed. `10 % 2 == 0` is true, so the `if` branch is executed and `n` becomes 5. By the end of the trace we have:

n	output printed so far
--	-----
3	3,
10	3, 10,
5	3, 10, 5,
16	3, 10, 5, 16,
8	3, 10, 5, 16, 8,
4	3, 10, 5, 16, 8, 4,
2	3, 10, 5, 16, 8, 4, 2,
1	3, 10, 5, 16, 8, 4, 2, 1.

Tracing can be a bit tedious and error prone (that's why we get computers to do this stuff in the first place!), but it is an essential skill for a programmer to have. From this trace we can learn a lot about the way our code works. We can observe that as soon as `n` becomes a power of 2, for example, the program will require  $\log_2(n)$  executions of the loop body to complete. We can also see that the final 1 will not be printed as output within the body of the loop, which is why we put the special `print` function at the end.

### 3.3.7 Counting digits

The following snippet counts the number of decimal digits in a positive integer:

```
1 n = 3029
2 count = 0
3 while n != 0:
4     count = count + 1
5     n = n // 10
6 print(count)
```

Trace the execution to convince yourself that it works.

This snippet demonstrates an important pattern of computation called a **counter**. The variable `count` is initialized to 0 and then incremented each time the loop body is executed. When the loop exits, `count` contains the result — the total number of times the loop body was executed, which is the same as the number of digits.

If we wanted to only count digits that are either 0 or 5, adding a conditional before incrementing the counter will do the trick:

```
1 n = 2574301453
2 count = 0
3 while n > 0:
4     digit = n % 10
5     if digit == 0 or digit == 5:
6         count = count + 1
7     n = n // 10
8 print(count)
```

Notice, however, that if `n = 0` this snippet will not print 1 as answer. Explain why. Do you think this is a bug in the code, or a bug in the specifications, or our expectations?

### 3.3.8 Help and meta-notation

Python comes with extensive documentation for all its built-in functions, and its libraries. Different systems have different ways of accessing this help. See for example <https://docs.python.org/3/library/stdtypes.html#typeseq-range>

Notice the square brackets in the description of the arguments. These are examples of **meta-notation** — notation that describes Python syntax, but is not part of it. The square brackets in this documentation mean that the argument is *optional* — the programmer can omit it. So what this first line of help tells us is that `range` must always have a `stop` argument, but it may have an optional `start` argument (which must be followed by a comma if it is present), and it can also have an optional `step` argument, preceded by a comma if it is present.

The examples from help show that `range` can have either 1, 2 or 3 arguments. The list can start at any starting value, and go up or down in increments other than 1. The documentation here also says that the arguments must be integers.

Other meta-notation you'll frequently encounter is the use of bold and italics. The bold means that these are tokens — keywords or symbols — typed into your Python code exactly as they are, whereas the italic terms stand for “something of this type”. So the syntax description

**for** *variable* **in** *list* :

means you can substitute any legal variable and any legal list when you write your Python code.

This (simplified) description of the `print` function, shows another example of meta-notation in which the ellipses (`...`) mean that you can have as many objects as you like (even zero), separated by commas:

**print**( [*object*, ... ] )

Meta-notation gives us a concise and powerful way to describe the *pattern* of some syntax or feature.



### 3.3.9 Tables

One of the things loops are good for is generating tables. Before computers were readily available, people had to calculate logarithms, sines and cosines, and other mathematical functions by hand. To make that easier, mathematics books contained long tables listing the values of these functions. Creating the tables was slow and boring, and they tended to be full of errors.

When computers appeared on the scene, one of the initial reactions was, “*This is great! We can use the computers to generate the tables, so there will be no errors.*” That turned out to be true (mostly) but shortsighted. Soon thereafter, computers and calculators were so pervasive that the tables became obsolete.

Well, almost. For some operations, computers use tables of values to get an approximate answer and then perform computations to improve the approximation. In some cases, there have been errors in the underlying tables, most famously in the table the Intel Pentium processor chip used to perform floating-point division.

Although a log table is not as useful as it once was, it still makes a good example of iteration. The following program outputs a sequence of values in the left column and 2 raised to the power of that value in the right column:

```
1 for x in range(13):    # Generate numbers 0 to 12
2     print(x, "\t", 2**x)
```

The string “\t” represents a **tab character**. The backslash character in “\t” indicates the beginning of an **escape sequence**. Escape sequences are used to represent invisible characters like tabs and newlines. The sequence \n represents a **newline**.

An escape sequence can appear anywhere in a string; in this example, the tab escape sequence is the only thing in the string. How do you think you represent a backslash in a string?

As characters and strings are displayed on the screen, an invisible marker called the **cursor** keeps track of where the next character will go. After a `print` function, the cursor normally goes to the beginning of the next line.

The tab character shifts the cursor to the right until it reaches one of the tab stops. Tabs are useful for making columns of text line up, as in the output of the previous program:

0	1
1	2
2	4
3	8
4	16
5	32
6	64
7	128
8	256
9	512
10	1024
11	2048
12	4096

Because of the tab characters between the columns, the position of the second column does not depend on the number of digits in the first column.

### 3.3.10 Two-dimensional tables

A two-dimensional table is a table where you read the value at the intersection of a row and a column. A multiplication table is a good example. Let’s say you want to print a multiplication table for the values from 1 to 6.

A good way to start is to write a loop that prints the multiples of 2, all on one line:

```
1 for i in range(1, 7):
2     print(2 * i, end=" ")
3 print()
```

Here we've used the `range` function, but made it start its sequence at 1. As the loop executes, the value of `i` changes from 1 to 6. When all the elements of the range have been assigned to `i`, the loop terminates. Each time through the loop, it displays the value of `2 * i`, followed by three spaces.

Again, the extra `end=" "` argument in the `print` function suppresses the newline, and uses three spaces instead. After the loop completes, the call to `print` at line 3 finishes the current line, and starts a new line.

The output of the program is:

```
2      4      6      8     10     12
```

So far, so good. The next step is to **encapsulate** and **generalize**. We will continue this topic in the next chapter.

### 3.3.11 The `break` statement

The **`break`** statement is used to immediately leave the body of its loop. The next statement to be executed is the first one after the body:

```
1 for i in [12, 16, 17, 24, 29]:
2     if i % 2 == 1: # If the number is odd
3         break    # ... immediately exit the loop
4     print(i)
5 print("done")
```

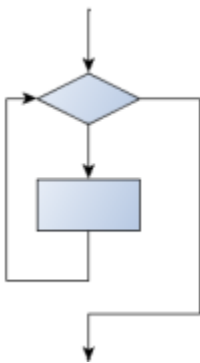
This prints:

```
12
16
done
```

---

#### The pre-test loop — standard loop behaviour

`for` and `while` loops do their tests at the start, before executing any part of the body. They're called **pre-test** loops, because the test happens before (pre) the body. `break` and `return` (discussed later) are our tools for adapting this standard behaviour.



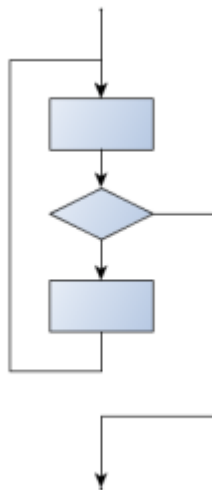
### 3.3.12 Other flavours of loops

Sometimes we'd like to have the **middle-test** loop with the exit test in the middle of the body, rather than at the beginning or at the end. Or a **post-test** loop that puts its exit test as the last thing in the body. Other languages have different syntax and keywords for these different flavours, but Python just uses a combination of `while` and `if` `<CONDITION>`: `break` to get the job done.

A typical example is a problem where the user has to input numbers to be summed. To indicate that there are no more inputs, the user enters a special value, often the value `-1`, or the empty string. This needs a middle-exit loop pattern: input the next number, then test whether to exit, or else process the number:

---

The middle-test loop flowchart



---

```
1 total = 0
2 while True:
3     response = input("Enter the next number. (Leave blank to end)")
4     if response == "" or response == "-1":
5         break
6     total += int(response)
7 print("The total of the numbers you entered is ", total)
```

Convince yourself that this fits the middle-exit loop flowchart: line 3 does some useful work, lines 4 and 5 can exit the loop, and if they don't line 6 does more useful work before the next iteration starts.

The `while bool-expr:` uses the Boolean expression to determine whether to iterate again. `True` is a trivial Boolean expression, so `while True:` means *always do the loop body again*. This is a language *idiom* — a convention that most programmers will recognize immediately. Since the expression on line 2 will never terminate the loop, (it is a dummy test) the programmer must arrange to break (or return) out of the loop body elsewhere, in some other way (i.e. in lines 4 and 5 in this sample). A clever compiler or interpreter will understand that line 2 is a fake test that must always succeed, so it won't even generate a test, and our flowchart never even put the diamond-shape dummy test box at the top of the loop!

Similarly, by just moving the `if condition: break` to the end of the loop body we create a pattern for a post-test loop. Post-test loops are used when you want to be sure that the loop body always executes at least once (because the first test only happens at the end of the execution of the first loop body). This is useful, for example, if we want to play an interactive game against the user — we always want to play at least one game:

```
1 while True:
2     play_the_game_once()
3     response = input("Play again? (yes or no)")
4     if response != "yes":
5         break
6 print("Goodbye!")
```

---

**Hint: Think about where you want the exit test to happen**

Once you’ve recognized that you need a loop to repeat something, think about its terminating condition — when will I want to stop iterating? Then figure out whether you need to do the test before starting the first (and every other) iteration, or at the end of the first (and every other) iteration, or perhaps in the middle of each iteration. Interactive programs that require input from the user or read from files often need to exit their loops in the middle or at the end of an iteration, when it becomes clear that there is no more data to process, or the user doesn’t want to play our game anymore.

---

### 3.3.13 An example

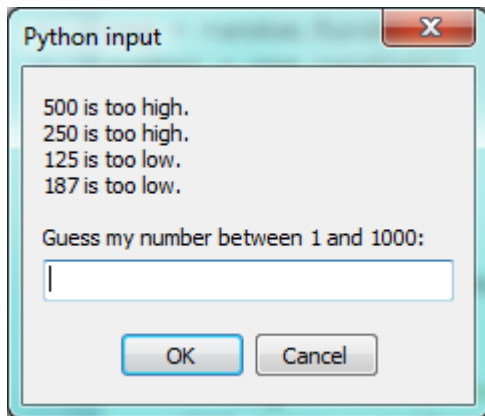
The following program implements a simple guessing game:

```
1 import random                                # We cover random numbers in the
2 rng = random.Random()                        # modules chapter, so peek ahead if you
3 ↪ want. "rng" stands for "random number generator".
4 number = rng.randrange(1, 1000) # Get random number between [1 and 1000).
5
6 guesses = 0
7 message = ""
8
9 while True:
10     guess = int(input(message + "\nGuess my number between 1 and 1000: "))
11     guesses += 1
12     if guess > number:
13         message += str(guess) + " is too high.\n"
14     elif guess < number:
15         message += str(guess) + " is too low.\n"
16     else:
17         break
18 input("\n\nGreat, you got it in "+str(guesses)+" guesses!\n\n")
```

This program makes use of the mathematical law of **trichotomy** (given real numbers  $a$  and  $b$ , exactly one of these three must be true:  $a > b$ ,  $a < b$ , or  $a == b$ ).

At line 18 there is a call to the input function, but we don’t do anything with the result, not even assign it to a variable. This is legal in Python. Here it has the effect of popping up the input dialog window and waiting for the user to respond before the program terminates. Programmers often use the trick of doing some extra input at the end of a script, just to keep the window open.

Also notice the use of the `message` variable, initially an empty string, on lines 6, 12 and 14. Each time through the loop we extend the message being displayed: this allows us to display the program’s feedback right at the same place as we’re asking for the next guess.



### 3.3.14 The `continue` statement

This is a control flow statement that causes the program to immediately skip the processing of the rest of the body of the loop, *for the current iteration*. But the loop still carries on running for its remaining iterations:

```
1 for i in [12, 16, 17, 24, 29, 30]:
2     if i % 2 == 1:      # If the number is odd
3         continue      # Don't process it
4     print(i)
5 print("done")
```

This prints:

```
12
16
24
30
done
```

### 3.3.15 Paired Data

We've already seen lists of names and lists of numbers in Python. We're going to peek ahead in the textbook a little, and show a more advanced way of representing our data. Making a pair of things in Python is as simple as putting them into parentheses, like this:

```
1 year_born = ("Paris Hilton", 1981)
```

We can put many pairs into a list of pairs:

```
1 celebs = [("Brad Pitt", 1963), ("Jack Nicholson", 1937),
2           ("Justin Bieber", 1994)]
```

Here is a quick sample of things we can do with structured data like this. First, print all the celebs:

```
1 print(celebs)
2 print(len(celebs))
```

```
[("Brad Pitt", 1963), ("Jack Nicholson", 1937), ("Justin Bieber", 1994)]
3
```

Notice that the `celebs` list has just 3 elements, each of them pairs.

Now we print the names of those celebrities born before 1980:

```
1 for name, year in celebs:
2     if year < 1980:
3         print(name)
```

```
Brad Pitt
Jack Nicholson
```

This demonstrates something we have not seen yet in the `for` loop: instead of using a single loop control variable, we've used a pair of variable names, `(name, year)`, instead. The loop is executed three times — once for each pair in the list, and on each iteration both the variables are assigned values from the pair of data that is being handled.

### 3.3.16 Nested Loops for Nested Data

Now we'll come up with an even more adventurous list of structured data. In this case, we have a list of students. Each student has a name which is paired up with another list of subjects that they are enrolled for:

```
1 students = [
2     ("John", ["CompSci", "Physics"]),
3     ("Vusi", ["Maths", "CompSci", "Stats"]),
4     ("Jess", ["CompSci", "Accounting", "Economics", "Management"]),
5     ("Sarah", ["InfSys", "Accounting", "Economics", "CommLaw"]),
6     ("Zuki", ["Sociology", "Economics", "Law", "Stats", "Music"])]
```

Here we've assigned a list of five elements to the variable `students`. Let's print out each student name, and the number of subjects they are enrolled for:

```
1 # Print all students with a count of their courses.
2 for name, subjects in students:
3     print(name, "takes", len(subjects), "courses")
```

Python agreeably responds with the following output:

```
John takes 2 courses
Vusi takes 3 courses
Jess takes 4 courses
Sarah takes 4 courses
Zuki takes 5 courses
```

Now we'd like to ask how many students are taking `CompSci`. This needs a counter, and for each student we need a second loop that tests each of the subjects in turn:

```
1 # Count how many students are taking CompSci
2 counter = 0
3 for name, subjects in students:
4     for s in subjects:           # A nested loop!
5         if s == "CompSci":
6             counter += 1
7
8 print("The number of students taking CompSci is", counter)
```

```
The number of students taking CompSci is 3
```

A more concise of doing this would be the following:

```
1 counter = 0
2 for name, subjects in students:
3     if "CompSci" in subjects:
4         counter += 1
```

You should set up a list of your own data that interests you — perhaps a list of your CDs, each containing a list of song titles on the CD, or a list of movie titles, each with a list of movie stars who acted in the movie. You could then ask questions like “Which movies starred Angelina Jolie?”

### 3.3.17 Newton’s method for finding square roots

Loops are often used in programs that compute numerical results by starting with an approximate answer and iteratively improving it.

For example, before we had calculators or computers, people needed to calculate square roots manually. Newton used a particularly good method (there is some evidence that this method was known many years before). Suppose that you want to know the square root of  $n$ . If you start with almost any approximation, you can compute a better approximation (closer to the actual answer) with the following formula:

```
1 better = (approximation + n/approximation)/2
```

Repeat this calculation a few times using your calculator. Can you see why each iteration brings your estimate a little closer? One of the amazing properties of this particular algorithm is how quickly it converges to an accurate answer — a great advantage for doing it manually.

By using a loop and repeating this formula until the better approximation gets close enough to the previous one, we can write a function for computing the square root. (In fact, this is how your calculator finds square roots — it may have a slightly different formula and method, but it is also based on repeatedly improving its guesses.)

This is an example of an *indefinite* iteration problem: we cannot predict in advance how many times we’ll want to improve our guess — we just want to keep getting closer and closer. Our stopping condition for the loop will be when our old guess and our improved guess are “close enough” to each other.

Ideally, we’d like the old and new guess to be exactly equal to each other when we stop. But exact equality is a tricky notion in computer arithmetic when real numbers are involved. Because real numbers are not represented absolutely accurately (after all, a number like pi or the square root of two has an infinite number of decimal places because it is irrational), we need to formulate the stopping test for the loop by asking “is  $a$  close enough to  $b$ ”? This stopping condition can be coded like this:

```
1 threshold = 0.001
2 if abs(a-b) < threshold: # Make this smaller for better accuracy
3     break
```

Notice that we take the absolute value of the difference between  $a$  and  $b$ !

This problem is also a good example of when a middle-exit loop is appropriate:

```
1 n = 8
2 threshold = 0.001
3 approximation = n/2 # Start with some or other guess at the answer
4 while True:
5     better = (approximation + n/approximation)/2
6     if abs(approximation - better) < threshold:
7         print(better)
```

(continues on next page)

(continued from previous page)

```
8         break
9     approximation = better
```

See if you can improve the approximations by changing the stopping condition. Also, step through the algorithm (perhaps by hand, using your calculator) to see how many iterations were needed before it achieved this level of accuracy for `sqrt(25)`.

### 3.3.18 Algorithms

Newton’s method is an example of an **algorithm**: it is a mechanical process for solving a category of problems (in this case, computing square roots).

Some kinds of knowledge are not algorithmic. For example, learning dates from history or your multiplication tables involves memorization of specific solutions.

But the techniques you learned for addition with carrying, subtraction with borrowing, and long division are all algorithms. Or if you are an avid Sudoku puzzle solver, you might have some specific set of steps that you always follow.

One of the characteristics of algorithms is that they do not require any intelligence to carry out. They are mechanical processes in which each step follows from the last according to a simple set of rules. And they’re designed to solve a general class or category of problems, not just a single problem.

Understanding that hard problems can be solved by step-by-step algorithmic processes (and having technology to execute these algorithms for us) is one of the major breakthroughs that has had enormous benefits. So while the execution of the algorithm may be boring and may require no intelligence, algorithmic or computational thinking — i.e. using algorithms and automation as the basis for approaching problems — is rapidly transforming our society. Some claim that this shift towards algorithmic thinking and processes is going to have even more impact on our society than the invention of the printing press. And the process of designing algorithms is interesting, intellectually challenging, and a central part of what we call programming.

Some of the things that people do naturally, without difficulty or conscious thought, are the hardest to express algorithmically. Understanding natural language is a good example. We all do it, but so far no one has been able to explain *how* we do it, at least not in the form of a step-by-step mechanical algorithm.

## 3.4 Some Tips, Tricks, and Common Errors

These are small summaries of ideas, tips, and commonly seen errors that might be helpful to those beginning Python.

### 3.4.1 Problems with logic and flow of control

We often want to know if some condition holds for any item in a list, e.g. “does the list have any odd numbers?” This is a common mistake:

```
1 numbers = [10, 5, 24, 8, 6]
2
3 # Buggy version
4 for number in numbers:
5     if number % 2 == 1:
6         print(True)
7         break
8     else:
```

(continues on next page)



(continued from previous page)

```
9     print(False)
10    break
```

Can we spot two problems here? As soon as we execute a `break`, we'll leave the loop. So the logic of saying "If I find an odd number I can return `True`" is fine. However, we cannot return `False` after only looking at one item — we can only return `False` if we've been through all the items, and none of them are odd. So line 10 should not be there, and lines 8 and 9 have to be outside the loop. Here is a corrected version:

```
1 numbers = [10, 5, 24, 8, 6]
2 for number in numbers:
3     if number % 2 == 1:
4         print(True)
5         break
6 else:
7     print(False)
```

We'll see This "eureka", or "short-circuit" style of breaking from a loop as soon as we are certain what the outcome will be again later.

---

**Note that this uses a `for ... else` construct.**

The `else` clause is executed when a loop has looped without encountering any `break` statements. This is ideal for our case here. Also note that the `else` is not, in this case, related to the `if` statement that occurs inside the loop.

---

It is preferred over this one, which also works correctly:

```
1 numbers = [10, 5, 24, 8, 6]
2 count = 0
3 for number in numbers:
4     if number % 2 == 1:
5         count += 1      # Count the odd numbers
6 if count > 0:
7     print(True)
8 else:
9     print(False)
```

The performance disadvantage of this one is that it traverses the whole list, even if it knows the outcome very early on.

---

**Tip: Think about the return conditions of the loop**

Do I need to look at all elements in all cases? Can I shortcut and take an early exit? Under what conditions? When will I have to examine all the items in the list?

---

The code in lines 6-9 can also be tightened up. The expression `count > 0` itself represents a Boolean value, either `True` or `False` (we can say it 'evaluates' to either `True` or `False`). That `True/False` value can be used directly in the `print` statement. So we could cut out that code and simply have the following:

```
1 numbers = [10, 5, 24, 8, 6]
2 count = 0
3 for number in numbers:
4     if number % 2 == 1:
5         count += 1      # Count the odd numbers
6 print(count > 0)        # Aha! a programmer who understands that Boolean
7                          # expressions are not just used in if statements!
```

Although this code is tighter, it is not as nice as the one that did the short-circuit return as soon as the first odd number was found.

Even shorter:

```
1 numbers = [10, 5, 24, 8, 6]
2 count = 0
3 for number in numbers:
4     count += number % 2 == 1
5 print(count > 0)    # Aha! a programmer who understands that Boolean
6                     # expressions are not just used in if statements!
```

---

**Tip: Generalize your use of Booleans**

Programmers won't write `if is_prime(n) == True:` when they could say instead `if is_prime(n):`. Think more generally about Boolean values, not just in the context of `if` or `while` statements. Like arithmetic expressions, they have their own set of operators (`and`, `or`, `not`) and values (`True`, `False`) and can be assigned to variables, put into lists, etc. A good resource for improving your use of Booleans is [http://en.wikibooks.org/wiki/Non-Programmer%27s\\_Tutorial\\_for\\_Python\\_3/Boolean\\_Expressions](http://en.wikibooks.org/wiki/Non-Programmer%27s_Tutorial_for_Python_3/Boolean_Expressions)

---

Exercise time:

- How would we adapt this to print `True` if *all* the numbers are odd? Can you still use a short-circuit style?
- How would we adapt it to print `True` if at least three of the numbers are odd? Short-circuit the traversal when the third odd number is found — don't traverse the whole list unless we have to.

### 3.4.2 Looping and lists

Computers are useful because they can repeat computation, accurately and fast. So loops are going to be a central feature of almost all programs you encounter.

---

**Tip: Don't create unnecessary lists**

Lists are useful if you need to keep data for later computation. But if you don't need lists, it is probably better not to generate them.

---

Here are two functions that both generate ten million random numbers, and return the sum of the numbers. They both work.

```
1 import random
2 joe = random.Random()
3
4 # Version 1
5 # Build a list of random numbers, then sum them
6 numbers = []
7 for _ in range(10000000):
8     num = joe.randrange(1000)    # Generate one random number
9     numbers.append(num)          # Save it in our list, see the next_
    ↪chapter
10
11 tot = sum(numbers)
12 print(tot)
13
```

(continues on next page)

(continued from previous page)

```
14 # Version 2
15 # Sum the random numbers as we generate them
16 tot = 0
17 for _ in range(10000000):
18     num = joe.randrange(1000)
19     tot += num
20 print(tot)
```

What reasons are there for preferring the second version here? (Hint: open a tool like the Performance Monitor on your computer, and watch the memory usage. How big can you make the list before you get a fatal memory error in the first version?)

In a similar way, when working with files, we often have an option to read the whole file contents into a single string, or we can read one line at a time and process each line as we read it. Line-at-a-time is the more traditional and perhaps safer way to do things — you’ll be able to work comfortably no matter how large the file is. (And, of course, this mode of processing the files was essential in the old days when computer memories were much smaller.) But you may find whole-file-at-once is sometimes more convenient!

### 3.4.3 Glossary

**attribute** Some state or value that belongs to a particular object. For example, `tess` has a `color`.

**canvas** A surface within a window where drawing takes place.

**control flow** See *flow of execution*.

**for loop** A statement in Python for convenient repetition of statements in the *body* of the loop.

**loop body** Any number of statements nested inside a loop. The nesting is indicated by the fact that the statements are indented under the `for` loop statement.

**loop variable** A variable used as part of a `for` loop. It is assigned a different value on each iteration of the loop.

**instance** An object of a certain type, or class. `tess` and `alex` are different instances of the class `Turtle`.

**method** A function that is attached to an object. Invoking or activating the method causes the object to respond in some way, e.g. `forward` is the method when we say `tess.forward(100)`.

**invoke** An object has methods. We use the verb *invoke* to mean *activate the method*. Invoking a method is done by putting parentheses after the method name, with some possible arguments. So `tess.forward()` is an invocation of the `forward` method.

**module** A file containing Python definitions and statements intended for use in other Python programs. The contents of a module are made available to the other program by using the `import` statement.

**object** A “thing” to which a variable can refer. This could be a screen window, or one of the turtles we have created.

**range** A built-in function in Python for generating sequences of integers. It is especially useful when we need to write a `for` loop that executes a fixed number of times.

**terminating condition** A condition that occurs which causes a loop to stop repeating its body. In the `for` loops we saw in this chapter, the terminating condition has been when there are no more elements to assign to the loop variable.

**block** A group of consecutive statements with the same indentation.

**body** The block of statements in a compound statement that follows the header.

**Boolean algebra** Some rules for rearranging and reasoning about Boolean expressions.

**Boolean expression** An expression that is either true or false.

**Boolean value** There are exactly two Boolean values: `True` and `False`. Boolean values result when a Boolean expression is evaluated by the Python interpreter. They have type `bool`.

**branch** One of the possible paths of the flow of execution determined by conditional execution.

**chained conditional** A conditional branch with more than two possible flows of execution. In Python chained conditionals are written with `if ... elif ... else` statements.

**comparison operator** One of the six operators that compares two values: `==`, `!=`, `>`, `<`, `>=`, and `<=`.

**condition** The Boolean expression in a conditional statement that determines which branch is executed.

**conditional statement** A statement that controls the flow of execution depending on some condition. In Python the keywords `if`, `elif`, and `else` are used for conditional statements.

**logical operator** One of the operators that combines Boolean expressions: `and`, `or`, and `not`.

**nesting** One program structure within another, such as a conditional statement inside a branch of another conditional statement.

**prompt** A visual cue that tells the user that the system is ready to accept input data.

**truth table** A concise table of Boolean values that can describe the semantics of an operator.

**type conversion** An explicit function call that takes a value of one type and computes a corresponding value of another type.

**algorithm** A step-by-step process for solving a category of problems.

**body** The statements inside a loop.

**bump** Programmer slang. Synonym for increment.

**continue statement** A statement that causes the remainder of the current iteration of a loop to be skipped. The flow of execution goes back to the top of the loop, evaluates the condition, and if this is true the next iteration of the loop will begin.

**counter** A variable used to count something, usually initialized to zero and incremented in the body of a loop.

**cursor** An invisible marker that keeps track of where the next character will be printed.

**decrement** Decrease by 1.

**definite iteration** A loop where we have an upper bound on the number of times the body will be executed. Definite iteration is usually best coded as a `for` loop.

**escape sequence** An escape character, `\`, followed by one or more printable characters used to designate a nonprintable character.

**increment** Both as a noun and as a verb, increment means to increase by 1.

**infinite loop** A loop in which the terminating condition is never satisfied.

**indefinite iteration** A loop where we just need to keep going until some condition is met. A `while` statement is used for this case.

**initialization (of a variable)** To initialize a variable is to give it an initial value. Since in Python variables don't exist until they are assigned values, they are initialized when they are created. In other programming languages this is not the case, and variables can be created without being initialized, in which case they have either default or *garbage* values.

**iteration** Repeated execution of a set of programming statements.

**loop** The construct that allows us to repeatedly execute a statement or a group of statements until a terminating condition is satisfied.

**loop variable** A variable used as part of the terminating condition of a loop.

**meta-notation** Extra symbols or notation that helps describe other notation. Here we introduced square brackets, ellipses, italics, and bold as meta-notation to help describe optional, repeatable, substitutable and fixed parts of the Python syntax.

**middle-test loop** A loop that executes some of the body, then tests for the exit condition, and then may execute some more of the body. We don't have a special Python construct for this case, but can use `while` and `break` together.

**nested loop** A loop inside the body of another loop.

**newline** A special character that causes the cursor to move to the beginning of the next line.

**post-test loop** A loop that executes the body, then tests for the exit condition. We don't have a special Python construct for this, but can use `while` and `break` together.

**pre-test loop** A loop that tests before deciding whether to execute its body. `for` and `while` are both pre-test loops.

**tab** A special character that causes the cursor to move to the next tab stop on the current line.

**trichotomy** Given any real numbers  $a$  and  $b$ , exactly one of the following relations holds:  $a < b$ ,  $a > b$ , or  $a == b$ . Thus when you can establish that two of the relations are false, you can assume the remaining one is true.

**trace** To follow the flow of execution of a program by hand, recording the change of state of the variables and any output produced.

### 3.4.4 Exercises

1. Assume the days of the week are numbered 0,1,2,3,4,5,6 from Sunday to Saturday. Write a program that asks a day number, and prints the day name (a string).
2. You go on a wonderful holiday (perhaps to jail, if you don't like happy exercises) leaving on day number 3 (a Wednesday). You return home after 137 sleeps. Write a general version of the program which asks for the starting day number, and the length of your stay, and it will tell you the name of day of the week you will return on.
3. Give the logical opposites of these conditions
  - (a)  $a > b$
  - (b)  $a \geq b$
  - (c)  $a \geq 18$  and  $day == 3$
  - (d)  $a \geq 18$  and  $day != 3$
4. What do these expressions evaluate to?
  - (a)  $3 == 3$
  - (b)  $3 != 3$
  - (c)  $3 \geq 4$
  - (d) `not (3 < 4)`
5. Complete this truth table:

p	q	r	(not (p and q)) or r
F	F	F	?
F	F	T	?
F	T	F	?
F	T	T	?
T	F	F	?
T	F	T	?
T	T	F	?
T	T	T	?

6. Write a program which is given an exam mark, and it returns a string — the grade for that mark — according to this scheme:

Mark	Grade
$\geq 75$	First
[70-75)	Upper Second
[60-70)	Second
[50-60)	Third
[45-50)	F1 Supp
[40-45)	F2
$< 40$	F3

The square and round brackets denote closed and open intervals. A closed interval includes the number, and open interval excludes it. So 39.99999 gets grade F3, but 40 gets grade F2. Assume

```
numbers = [83, 75, 74.9, 70, 69.9, 65, 60, 59.9, 55, 50,
           49.9, 45, 44.9, 40, 39.9, 2, 0]
```

Test your code by printing the mark and the grade for all the elements in this list.

7. Write a program which, given the length of two sides of a right-angled triangle, returns the length of the hypotenuse. (Hint:  $x ** 0.5$  will return the square root.)
8. Write a program which, given the length of three sides of a triangle, will determine whether the triangle is right-angled. Assume that the third argument to the function is always the longest side. It will return `True` if the triangle is right-angled, or `False` otherwise.

Hint: Floating point arithmetic is not always exactly accurate, so it is not safe to test floating point numbers for equality. If a good programmer wants to know whether  $x$  is equal or close enough to  $y$ , they would probably code it up as:

```
threshold = 1e-7
if abs(x-y) < threshold:    # If x is approximately equal to y
    ...
```

9. Extend the above program so that the sides can be given to the function in any order.
10. If you're intrigued by why floating point arithmetic is sometimes inaccurate, on a piece of paper, divide 10 by 3 and write down the decimal result. You'll find it does not terminate, so you'll need an infinitely long sheet of paper. The *representation* of numbers in computer memory or on your calculator has similar problems: memory is finite, and some digits may have to be discarded. So small inaccuracies creep in. Try this script:

```
1 import math
2 a = math.sqrt(2.0)
```

(continues on next page)

(continued from previous page)

```
3 print(a, a*a)
4 print(a*a == 2.0)
```

1. Write a program that prints We like Python's turtles! 1000 times.
2. **Write a program that uses a for loop to print**  
  
One of the months of the year is January  
One of the months of the year is February  
...
  - (a) Write a loop that prints each of the numbers on a new line.
  - (b) Write a loop that prints each number and its square on a new line.
  - (c) Write a loop that adds all the numbers from the list into a variable called *total*. You should set the *total* variable to have the value 0 before you start adding them up, and print the value in *total* after the loop has completed.
  - (d) Print the product of all the numbers in the list. (product means all multiplied together)
5. Use for loops to make a turtle draw these regular polygons (regular means all sides the same lengths, all angles the same):
  - An equilateral triangle
  - A square
  - A hexagon (six sides)
  - An octagon (eight sides)
6. A drunk pirate makes a random turn and then takes 100 steps forward, makes another random turn, takes another 100 steps, turns another random amount, etc. A social science student records the angle of each turn before the next 100 steps are taken. Her experimental data is [160, -43, 270, -97, -43, 200, -940, 17, -86]. (Positive angles are counter-clockwise.) Use a turtle to draw the path taken by our drunk friend.
7. Enhance your program above to also tell us what the drunk pirate's heading is after he has finished stumbling around. (Assume he begins at heading 0).
8. If you were going to draw a regular polygon with 18 sides, what angle would you need to turn the turtle at each corner?
9. At the interactive prompt, anticipate what each of the following lines will do, and then record what happens. Score yourself, giving yourself one point for each one you anticipate correctly:

```
>>> import turtle
>>> window = turtle.Screen()
>>> tess = turtle.Turtle()
>>> tess.right(90)
>>> tess.left(3600)
>>> tess.right(-90)
>>> tess.speed(10)
>>> tess.left(3600)
>>> tess.speed(0)
>>> tess.left(3645)
>>> tess.forward(-100)
```

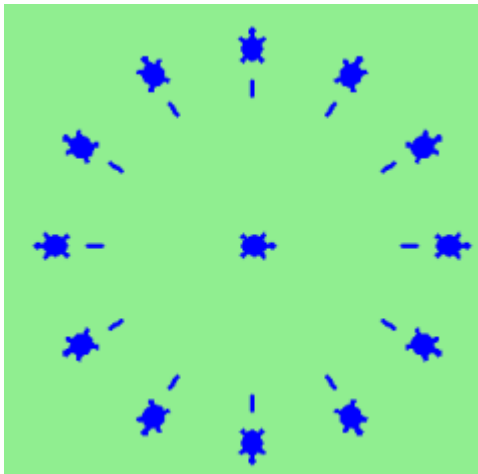
10. Write a program to draw a shape like this:



Hints:

- Try this on a piece of paper, moving and turning your cellphone as if it was a turtle. Watch how many complete rotations your cellphone makes before you complete the star. Since each full rotation is 360 degrees, you can figure out the total number of degrees that your phone was rotated through. If you divide that by 5, because there are five points to the star, you'll know how many degrees to turn the turtle at each point.
- You can hide a turtle behind its invisibility cloak if you don't want it shown. It will still draw its lines if its pen is down. The method is invoked as `tess.hideturtle()`. To make the turtle visible again, use `tess.showturtle()`.

11. Write a program to draw a face of a clock that looks something like this:



12. Create a turtle, and assign it to a variable. When you ask for its type, what do you get?

This chapter showed us how to sum a list of items, and how to count items. The counting example also had an `if` statement that let us only count some selected items. we have `break` to exit a loop, and `continue` to abandon the current iteration of the loop without ending the loop.

Composition of list traversal, summing, counting, testing conditions and early exit is a rich collection of building blocks that can be combined in powerful ways to create many functions that are all slightly different.

The first six questions are typical functions you should be able to write using only these building blocks.

1. Write a program to count how many odd numbers are in a list.
2. Sum up all the even numbers in a list.
3. Sum up all the negative numbers in a list.
4. Count how many words in a list have length 5.
5. Sum all the elements in a list up to but not including the first even number. (What if there is no even number?)



6. Count how many words occur in a list up to and including the first occurrence of the word “sam”. (What if “sam” does not occur?)
7. Add a print function to Newton’s `sqrt` algorithm that prints out `better` each time it is calculated. Call your modified program with 25 as an argument and record the results.
8. Write a program that prints out the first `n` triangular numbers. A call to with `n = 5` would produce the following output:

```
1      1
2      3
3      6
4     10
5     15
```

(*hint: use a web search to find out what a triangular number is.*)

9. Write a program which prints `True` when `n` is a *prime number* and `False` otherwise.
10. Revisit the drunk pirate problem. This time, the drunk pirate makes a turn, and then takes some steps forward, and repeats this. Our social science student now records *pairs* of data: the angle of each turn, and the number of steps taken after the turn. Her experimental data is [(160, 20), (-43, 10), (270, 8), (-43, 12)]. Use a turtle to draw the path taken by our drunk friend.
11. Many interesting shapes can be drawn by the turtle by giving a list of pairs like we did above, where the first item of the pair is the angle to turn, and the second item is the distance to move forward. Set up a list of pairs so that the turtle draws a house with a cross through the centre, as show here. This should be done without going over any of the lines / edges more than once, and without lifting your pen.



12. Recall the digit counting program. What will it print with `n = 0`? Modify it to print 1 for this case. Why does a call with `n = -24` result in an infinite loop? (*hint: `-1//10` evaluates to `-1`*) Modify `num_digits` so that it works correctly with any integer value.
13. Write a program that counts the number of even digits in `n`.
14. Write a program that computes the sum of the squares of the numbers in the list `numbers`. For example a call with, `numbers = [2, 3, 4]` should print 4+9+16 which is 29.

## 4.1 Functions

In Python, a **function** is a named sequence of statements that belong together. Their primary purpose is to help us organize programs into chunks that match how we think about the problem.

The syntax for a **function definition** is:

```
def <NAME>( <PARAMETERS> ) :  
    <STATEMENTS>
```

We can make up any names we want for the functions we create, except that we can't use a name that is a Python keyword, and the names must follow the rules for legal identifiers.

There can be any number of statements inside the function, but they have to be indented from the `def`. In the examples in this book, we will use the standard indentation of four spaces. Function definitions are the second of several **compound statements** we will see, all of which have the same pattern:

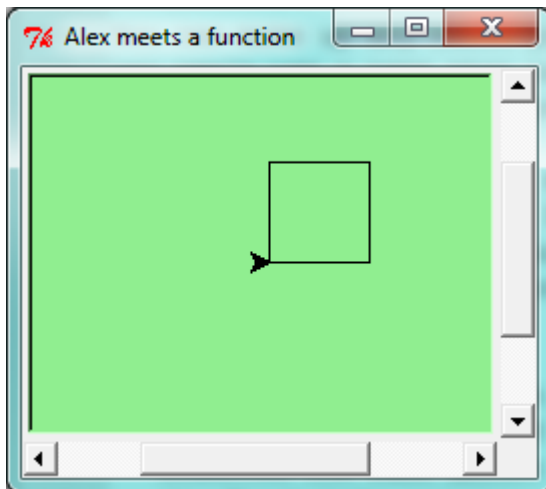
1. A header line which begins with a keyword and ends with a colon.
2. A **body** consisting of one or more Python statements, each indented the same amount — *the Python style guide recommends 4 spaces* — from the header line.

We've already seen the `for` loop which follows this pattern.

So looking again at the function definition, the keyword in the header is `def`, which is followed by the name of the function and some *parameters* enclosed in parentheses. The parameter list may be empty, or it may contain any number of parameters separated from one another by commas. In either case, the parentheses are required. The parameters specifies what information, if any, we have to provide in order to use the new function.

Suppose we're working with turtles, and a common operation we need is to draw squares. "Draw a square" is an *abstraction*, or a mental chunk, of a number of smaller steps. So let's write a function to capture the pattern of this "building block":

```
1 import turtle
2
3 def draw_square(animal, size):
4     """
5     Make animal draw a square with sides of length size.
6     """
7     for _ in range(4):
8         animal.forward(size)
9         animal.left(90)
10
11
12 window = turtle.Screen()           # Set up the window and its attributes
13 window.bgcolor("lightgreen")
14 window.title("Alex meets a function")
15
16 alex = turtle.Turtle()             # Create alex
17 draw_square(alex, 50)              # Call the function to draw the square
18 window.mainloop()
```



This function is named `draw_square`. It has two parameters: one to tell the function which turtle to move around, and the other to tell it the size of the square we want drawn. Make sure you know where the body of the function ends — it depends on the indentation, and the blank lines don't count for this purpose!

---

### Docstrings for documentation

If the first thing after the function header is a string, it is treated as a **docstring** and gets special treatment in Python and in some programming tools.

Docstrings are the key way to document our functions in Python and the documentation part is important. Because whoever calls our function shouldn't have to need to know what is going on in the function or how it works; they just need to know what arguments our function takes, what it does, and what the expected result is. Enough to be able to use the function without having to look underneath. This goes back to the concept of abstraction of which we'll talk more about.

Docstrings are usually formed using triple-quoted strings as they allow us to easily expand the docstring later on should we want to write more than a one-liner.

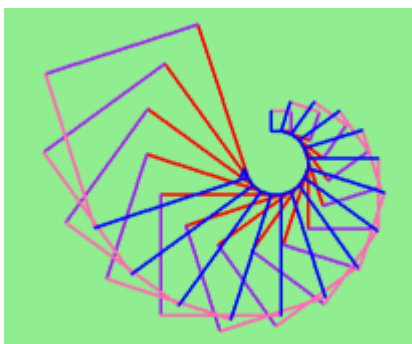
Just to differentiate from comments, a string at the start of a function (a docstring) is retrievable by Python tools *at runtime*. By contrast, comments are completely eliminated when the program is parsed.

---

Defining the function just tells Python *how* to do a particular task, not to *perform* it. In order to execute a function we need to make a **function call**. We've already seen how to call some built-in functions like **print**, **range** and **int**. Function calls contain the name of the function being executed followed by a list of values, called *arguments*, which are assigned to the parameters in the function definition. So in the second last line of the program, we call the function, and pass `alex` as the turtle to be manipulated, and 50 as the size of the square we want. While the function is executing, then, the variable `size` refers to the value 50, and the variable `animal` refers to the same turtle instance that the variable `alex` refers to. We called it `animal` to signify that there is no meaning to the name you give a function argument.

Once we've defined a function, we can call it as often as we like, and its statements will be executed each time we call it. And we could use it to get any of our turtles to draw a square. In the next example, we've changed the `draw_square` function a little, and we get `tess` to draw 15 squares, with some variations.

```
1  import turtle
2
3  def draw_multicolor_square(animal, size):
4      """Make animal draw a multi-color square of given size."""
5      for color in ["red", "purple", "hotpink", "blue"]:
6          animal.color(color)
7          animal.forward(size)
8          animal.left(90)
9
10     window = turtle.Screen()           # Set up the window and its attributes
11     window.bgcolor("lightgreen")
12
13     tess = turtle.Turtle()              # Create tess and set some attributes
14     tess.pensize(3)
15
16     size = 20                           # Size of the smallest square
17     for _ in range(15):
18         draw_multicolor_square(tess, size)
19         size += 10                       # Increase the size for next time
20         tess.forward(10)                 # Move tess along a little
21         tess.right(18)                   # and give her some turn
22
23     window.mainloop()
```



## 4.2 Functions can call other functions

Let's assume now we want a function to draw a rectangle. We need to be able to call the function with different arguments for width and height. And, unlike the case of the square, we cannot repeat the same thing 4 times, because the four sides are not equal.

So we eventually come up with this rather nice code that can draw a rectangle.

```
1 def draw_rectangle(animal, width, height):
2     """Get animal to draw a rectangle of given width and height."""
3     for _ in range(2):
4         animal.forward(width)
5         animal.left(90)
6         animal.forward(height)
7         animal.left(90)
```

*Thinking like a scientist* involves looking for patterns and relationships. In the code above, we've done that to some extent. We did not just draw four sides. Instead, we spotted that we could draw the rectangle as two halves, and used a loop to repeat that pattern twice.

But now we might spot that a square is a special kind of rectangle. We already have a function that draws a rectangle, so we can use that to draw our square.

```
1 def draw_square(animal, size):           # A new version of draw_square
2     draw_rectangle(animal, size, size)
```

There are some points worth noting here:

- Functions can call other functions.
- Rewriting `draw_square` like this captures the relationship that we've spotted between squares and rectangles.
- A caller of this function might say `draw_square(tess, 50)`. The parameters of this function, `animal` and `size`, are assigned the values of the `tess` object, and the int 50 respectively.
- In the body of the function they are just like any other variable.
- When the call is made to `draw_rectangle`, the values in variables `animal` and `size` are fetched first, then the call happens. So as we enter the top of function `draw_rectangle`, its variable `animal` is assigned the `tess` object, and `width` and `height` in that function are both given the value 50.

So far, it may not be clear why it is worth the trouble to create all of these new functions. Actually, there are a lot of reasons, but this example demonstrates two:

1. Creating a new function gives us an opportunity to name a group of statements. Functions can simplify a program by hiding a complex computation behind a single command. The function (including its name) can capture our mental chunking, or *abstraction*, of the problem.
2. Creating a new function can make a program smaller by eliminating repetitive code.

As we might expect, we have to create a function before we can execute it. In other words, the function definition has to be executed before the function is called.

## 4.3 Flow of execution

In order to ensure that a function is defined before its first use, we have to know the order in which statements are executed, which is called the **flow of execution**.

Execution always begins at the first statement of the program. Statements are executed one at a time, in order from top to bottom.

Function definitions do not alter the flow of execution of the program, but remember that statements inside the function are not executed until the function is called. Although it is not common, we can define one function inside another. In this case, the inner definition isn't executed until the outer function is called.

Function calls are like a detour in the flow of execution. Instead of going to the next statement, the flow jumps to the first line of the called function, executes all the statements there, and then comes back to pick up where it left off.

That sounds simple enough, until we remember that one function can call another. While in the middle of one function, the program might have to execute the statements in another function. But while executing that new function, the program might have to execute yet another function!

Fortunately, Python is adept at keeping track of where it is, so each time a function completes, the program picks up where it left off in the function that called it. When it gets to the end of the program, it terminates.

What's the moral of this sordid tale? When we read a program, don't read from top to bottom. Instead, follow the flow of execution.

As a simple example, let's consider the following program:

```
1 import turtle
2
3 def draw_square(animal, size):
4     for _ in range(4):
5         animal.forward(size)
6         animal.left(90)
7
8 window = turtle.Screen()           # Set up the window and its attributes
9
10 tess = turtle.Turtle()            # Create tess and set some attributes
11
12 draw_square(tess, 50)
13
14 window.mainloop()
```

The Python interpreter reads this script line by line. At the first line the `turtle` module is imported. We then define `draw_square`, which contains the instructions for a given turtle to draw a square. However, nothing happens *yet*. We then go on to define a window, and our charming turtle `tess`. The next line calls `draw_square`, asking `tess` to draw a square with sides of length 50. Finally, `window.mainloop()` actually runs these executions, and you will see `tess` draw a square on the screen.

Being able to trace your program is a valuable skill for a programmer.

## 4.4 Functions that require arguments

Most functions require arguments: the arguments provide for generalization. For example, if we want to find the absolute value of a number, we have to indicate what the number is. Python has a built-in function for computing the absolute value:

```
>>> abs(5)
5
>>> abs(-5)
5
```

In this example, the arguments to the `abs` function are 5 and -5.

Some functions take more than one argument. For example the built-in function `pow` takes two arguments, the base and the exponent. Inside the function, the values that are passed get assigned to variables called **parameters**.

```
>>> pow(2, 3)
8
>>> pow(7, 4)
2401
```

Another built-in function that takes more than one argument is `max`.

```
>>> max(7, 11)
11
>>> max(4, 1, 17, 2, 12)
17
>>> max(3 * 11, 5**3, 512 - 9, 1024**0)
503
```

`max` can be passed any number of arguments, separated by commas, and will return the largest value passed. The arguments can be either simple values or expressions. In the last example, 503 is returned, since it is larger than 33, 125, and 1.

## 4.5 Functions that return values

All the functions in the previous section return values. Calling each of these functions generates a value, which we usually assign to a variable or use as part of an expression.

```
1 biggest = max(3, 7, 2, 5)
2 x = abs(3 - 11) + 10
```

So an important difference between these functions and one like `draw_square` is that `draw_square` was not executed because we wanted it to compute a value — on the contrary, we wrote `draw_square` because we wanted it to execute a sequence of steps that caused the turtle to draw.

A function that returns a value is called a **fruitful function** in this book. The opposite of a fruitful function is **void function** — one that is not executed for its resulting value, but is executed because it does something useful. (Languages like Java, C#, C and C++ use the term “void function”, other languages like Pascal call it a **procedure**.) Even though void functions are not executed for their resulting value, Python always wants to return something. So if the programmer doesn’t arrange to return a value, Python will automatically return the value `None`.

How do we write our own fruitful function? In the exercises at the end of chapter 2 we saw the standard formula for compound interest, which we’ll now write as a fruitful function:

$$A = P \left( 1 + \frac{r}{n} \right)^{nt}$$

Where,

- `P` = principal amount (initial investment)
- `r` = annual nominal interest rate (as a decimal)
- `n` = number of times the interest is compounded per year
- `t` = number of years

```
1 def final_amount(p, r, n, t):
2     """
3     Apply the compound interest formula to p
4     to produce the final amount.
5     """
6
7     a = p * (1 + r/n) ** (n*t)
8     return a          # This is new, and makes the function fruitful.
9
10 # now that we have the function above, let us call it.
11 toInvest = float(input("How much do you want to invest?"))
```

(continues on next page)

(continued from previous page)

```
12 fnl = final_amount(toInvest, 0.08, 12, 5)
13 print("At the end of the period you'll have", fnl)
```

- The **return** statement is followed an expression (*a* in this case). This expression will be evaluated and returned to the caller as the “fruit” of calling this function.
- We prompted the user for the principal amount. The type of `toInvest` is a string, but we need a number before we can work with it. Because it is money, and could have decimal places, we’ve used the `float` type converter function to parse the string and return a float.
- Notice how we entered the arguments for 8% interest, compounded 12 times per year, for 5 years.
- When we run this, we get the output

*At the end of the period you'll have 14898.457083*

This is a bit messy with all these decimal places, but remember that Python doesn’t understand that we’re working with money: it just does the calculation to the best of its ability, without rounding. Later we’ll see how to format the string that is printed in such a way that it does get nicely rounded to two decimal places before printing.

- The line `toInvest = float(input("How much do you want to invest?"))` also shows yet another example of *composition* — we can call a function like `float`, and its arguments can be the results of other function calls (like `input`) that we’ve called along the way.

Notice something else very important here. The name of the variable we pass as an argument — `toInvest` — has nothing to do with the name of the parameter — `p`. It is as if `p = toInvest` is executed when `final_amount` is called. It doesn’t matter what the value was named in the caller, in `final_amount` its name is `p`.

These short variable names are getting quite tricky, so perhaps we’d prefer one of these versions instead:

```
1 def final_amount_v2(principal_amount, nominal_percentage_rate,
2                     num_times_per_year, years):
3     a = principal_amount * (1 + nominal_percentage_rate /
4                           num_times_per_year) ** (num_times_per_year*years)
5     return a
6
7 def final_amount_v3(amount, rate, compounded, years):
8     a = amount * (1 + rate/compounded) ** (compounded*years)
9     return a
10
11 def final_amount_v4(amount, rate, compounded, years):
12     """
13     The a in final_amount_v3 was a useless assignment.
14     We might as well skip it.
15     """
16     return amount * (1 + rate/compounded) ** (compounded*years)
```

They all do the same thing. Use your judgement to write code that can be best understood by other humans! Short variable names should generally be avoided, unless when short variables make more sense. This happens in particular with mathematical equations, where it’s perfectly fine to use *x*, *y*, etc.

## 4.6 Variables and parameters are local

When we create a **local variable** inside a function, it only exists inside the function, and we cannot use it outside. For example, consider again this function:



```
1 def final_amount(p, r, n, t):
2     a = p * (1 + r/n) ** (n*t)
3     return a
```

If we try to use `a`, outside the function, we'll get an error:

```
>>> a
NameError: name 'a' is not defined
```

The variable `a` is local to `final_amount`, and is not visible outside the function.

Additionally, `a` only exists while the function is being executed — we call this its **lifetime**. When the execution of the function terminates, the local variables are destroyed.

Parameters are also local, and act like local variables. For example, the lifetimes of `p`, `r`, `n`, `t` begin when `final_amount` is called, and the lifetime ends when the function completes its execution.

So it is not possible for a function to set some local variable to a value, complete its execution, and then when it is called again next time, recover the local variable. Each call of the function creates new local variables, and their lifetimes expire when the function returns to the caller.

## 4.7 Turtles Revisited

Now that we have fruitful functions, we can focus our attention on reorganizing our code so that it fits more nicely into our mental chunks. This process of rearrangement is called **refactoring** the code.

Two things we're always going to want to do when working with turtles is to create the window for the turtle, and to create one or more turtles. We could write some functions to make these tasks easier in future:

```
1 import turtle
2
3 def make_window(color, title):
4     """
5     Set up the window with the given background color and title.
6     Returns the new window.
7     """
8     window = turtle.Screen()
9     window.bgcolor(color)
10    window.title(title)
11    return window
12
13
14 def make_turtle(color, size):
15     """
16     Set up a turtle with the given color and pensize.
17     Returns the new turtle.
18     """
19     animal = turtle.Turtle()
20     animal.color(color)
21     animal.pensize(size)
22     return animal
23
24
25 wn = make_window("lightgreen", "Tess and Alex dancing")
26 tess = make_turtle("hotpink", 5)
```

(continues on next page)

(continued from previous page)

```
27 alex = make_turtle("black", 1)
28 dave = make_turtle("yellow", 2)
```

The trick about refactoring code is to anticipate which things we are likely to want to change each time we call the function: these should become the parameters, or changeable parts, of the functions we write.

## 4.8 Glossary

**argument** A value provided to a function when the function is called. This value is assigned to the corresponding parameter in the function. The argument can be the result of an expression which may involve operators, operands and calls to other fruitful functions.

**body** The second part of a compound statement. The body consists of a sequence of statements all indented the same amount from the beginning of the header. The standard amount of indentation used within the Python community is 4 spaces.

**compound statement** A statement that consists of two parts:

1. header - which begins with a keyword determining the statement type, and ends with a colon.
2. body - containing one or more statements indented the same amount from the header.

The syntax of a compound statement looks like this:

```
keyword ... :
    statement
statement ...
```

**docstring** A special string that is attached to a function as its `__doc__` attribute. Tools like Spyder can use docstrings to provide documentation or hints for the programmer. When we get to modules, classes, and methods, we'll see that docstrings can also be used there.

**flow of execution** The order in which statements are executed during a program run.

**frame** A box in a stack diagram that represents a function call. It contains the local variables and parameters of the function.

**function** A named sequence of statements that performs some useful operation. Functions may or may not take parameters and may or may not produce a result.

**function call** A statement that executes a function. It consists of the name of the function followed by a list of arguments enclosed in parentheses.

**function composition** Using the output from one function call as the input to another.

**function definition** A statement that creates a new function, specifying its name, parameters, and the statements it executes.

**fruitful function** A function that returns a value when it is called.

**header line** The first part of a compound statement. A header line begins with a keyword and ends with a colon (:)

**import statement** A statement which permits functions and variables defined in another Python module to be brought into the environment of another script. To use the features of the turtle, we need to first import the turtle module.

**lifetime** Variables and objects have lifetimes — they are created at some point during program execution, and will be destroyed at some time.

**local variable** A variable defined inside a function. A local variable can only be used inside its function. Parameters of a function are also a special kind of local variable.

**parameter** A name used inside a function to refer to the value which was passed to it as an argument.

**refactor** A fancy word to describe reorganizing our program code, usually to make it more understandable. Typically, we have a program that is already working, then we go back to “tidy it up”. It often involves choosing better variable names, or spotting repeated patterns and moving that code into a function.

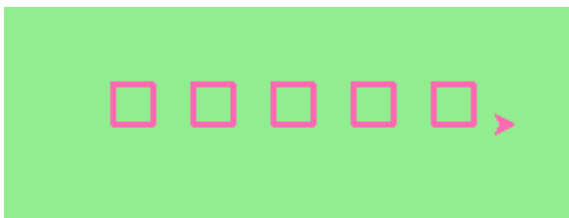
**stack diagram** A graphical representation of a stack of functions, their variables, and the values to which they refer.

**traceback** A list of the functions that are executing, printed when a runtime error occurs. A traceback is also commonly referred to as a *stack trace*, since it lists the functions in the order in which they are stored in the [runtime stack](#).

**void function** The opposite of a fruitful function: one that does not return a value. It is executed for the work it does, rather than for the value it returns.

## 4.9 Exercises

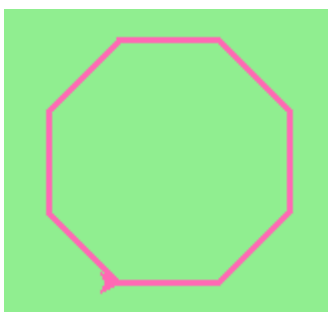
1. Write a void (non-fruitful) function to draw a square. Use it in a program to draw the image shown below. Assume each side is 20 units. (Hint: notice that the turtle has already moved away from the ending point of the last square when the program ends.)



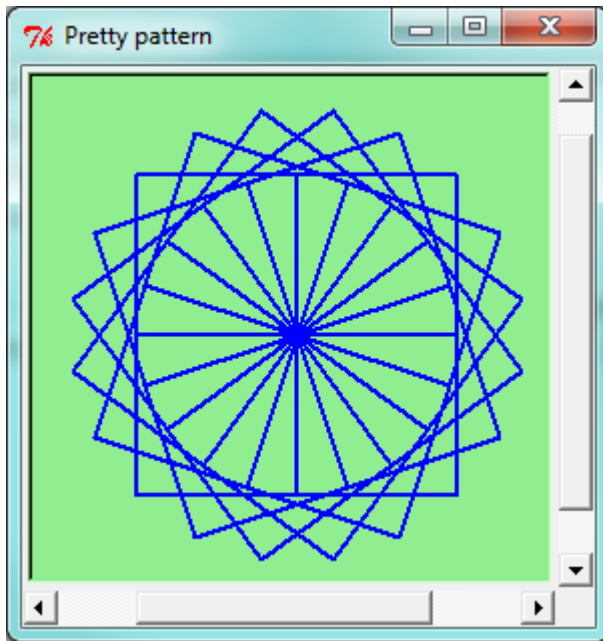
2. Write a program to draw this. Assume the innermost square is 20 units per side, and each successive square is 20 units bigger, per side, than the one inside it.



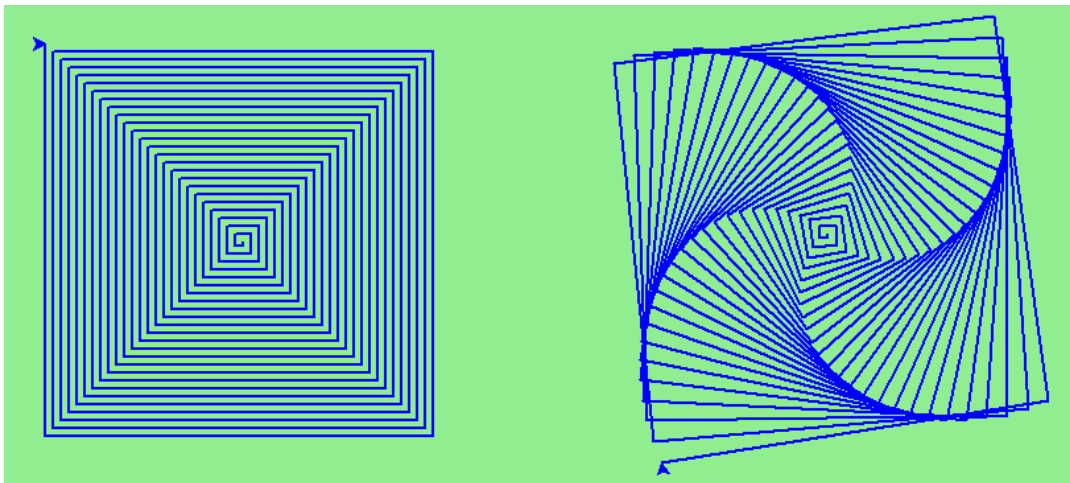
3. Write a void function `draw_poly(t, n, sz)` which makes a turtle draw a regular polygon. When called with `draw_poly(tess, 8, 50)`, it will draw a shape like this:



4. Draw this pretty pattern.



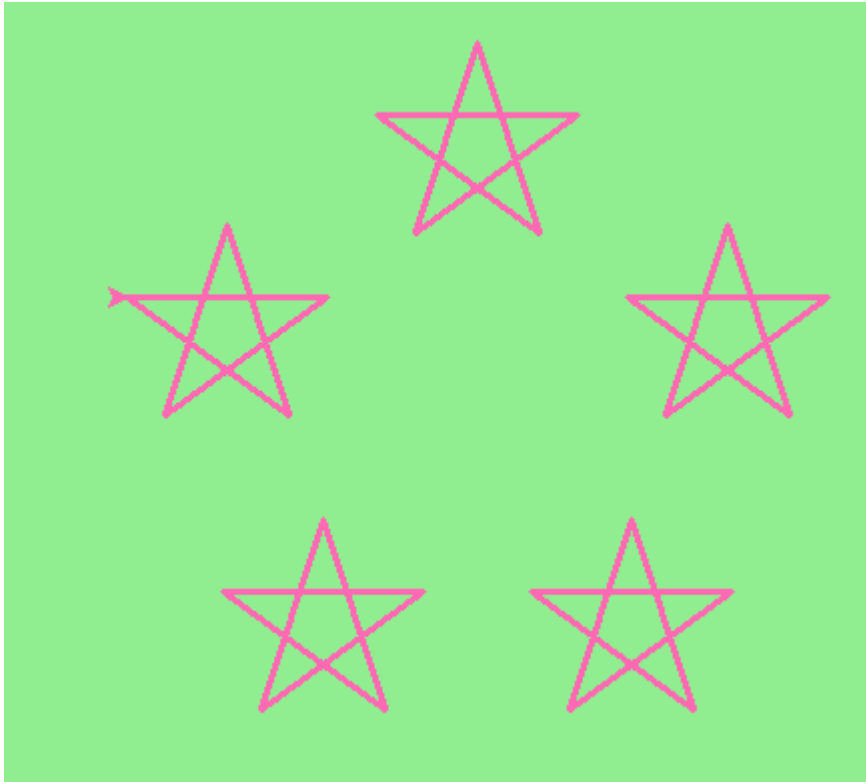
5. The two spirals in this picture differ only by the turn angle. Draw both.



6. Write a void function `draw_equitriangle(t, sz)` which calls `draw_poly` from the previous question to have its turtle draw an equilateral triangle.
7. Write a fruitful function `sum_to(n)` that returns the sum of all integer numbers up to and including `n`. So `sum_to(10)` would be `1+2+3...+10` which would return the value 55.
8. Write a function `area_of_circle(r)` which returns the area of a circle of radius `r`.
9. Write a void function to draw a star, where the length of each side is 100 units. (Hint: You should turn the turtle by 144 degrees at each point.)



10. Extend your program above. Draw five stars, but between each, pick up the pen, move forward by 350 units, turn right by 144, put the pen down, and draw the next star. You'll get something like this:



What would it look like if you didn't pick up the pen?

### 4.9.1 Fruitful functions

## 4.10 Return values

The built-in functions we have used, such as `abs`, `pow`, `int`, `max`, and `range`, have produced results. Calling each of these functions generates a value, which we usually assign to a variable or use as part of an expression.

```
1 biggest = max(3, 7, 2, 5)
2 x = abs(3 - 11) + 10
```

We also wrote our own function to return the final amount for a compound interest calculation.

In this chapter, we are going to write more functions that return values, which we will call *fruitful functions*, for want of a better name. The first example is `area`, which returns the area of a circle with the given radius:

```
1 def area(radius):
2     b = 3.14159 * radius**2
3     return b
```

We have seen the `return` statement before, but in a fruitful function the `return` statement includes a **return value**. This statement means: evaluate the return expression, and then return it immediately as the result (the fruit) of this function. The expression provided can be arbitrarily complicated, so we could have written this function like this:

```
1 def area(radius):
2     return 3.14159 * radius * radius
```

On the other hand, **temporary variables** like `b` above often make debugging easier.

Sometimes it is useful to have multiple return statements, one in each branch of a conditional. We have already seen the built-in `abs`, now we see how to write our own:

```
1 def absolute_value(x):
2     if x < 0:
3         return -x
4     else:
5         return x
```

Another way to write the above function is to leave out the `else` and just follow the `if` condition by the second `return` statement.

```
1 def absolute_value(x):
2     if x < 0:
3         return -x
4     return x
```

Think about this version and convince yourself it works the same as the first one.

Code that appears after a `return` statement, or any other place the flow of execution can never reach, is called **dead code**, or **unreachable code**.

In a fruitful function, it is a good idea to ensure that every possible path through the program hits a `return` statement. The following version of `absolute_value` fails to do this:

```
1 def bad_absolute_value(x):
2     if x < 0:
3         return -x
4     elif x > 0:
5         return x
```

This version is not correct because if `x` happens to be 0, neither condition is true, and the function ends without hitting a `return` statement. In this case, the return value is a special value called **None**:

```
>>> print(bad_absolute_value(0))
None
```

All Python functions return `None` whenever they do not return another value.

It is also possible to use a `return` statement in the middle of a `for` loop, in which case control immediately returns from the function. Let us assume that we want a function which looks through a list of words. It should return the first 2-letter word. If there is not one, it should return the empty string:

```
1 def find_first_2_letter_word(words):
2     for word in words:
```

(continues on next page)

(continued from previous page)

```
3     if len(word) == 2:
4         return word
5     return ""
```

```
>>> find_first_2_letter_word(["This", "is", "a", "dead", "parrot"])
'is'
>>> find_first_2_letter_word(["I", "like", "cheese"])
''
```

Single-step through this code and convince yourself that in the first test case that we've provided, the function returns while processing the second element in the list: it does not have to traverse the whole list.

## 4.11 Program development

At this point, you should be able to look at complete functions and tell what they do. Also, if you have been doing the exercises, you have written some small functions. As you write larger functions, you might start to have more difficulty, especially with runtime and semantic errors.

To deal with increasingly complex programs, we are going to suggest a technique called **incremental development**. The goal of incremental development is to avoid long debugging sessions by adding and testing only a small amount of code at a time.

As an example, suppose we want to find the distance between two points, given by the coordinates  $(x_1, y_1)$  and  $(x_2, y_2)$ . By the Pythagorean theorem, the distance is:

$$\text{distance} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

The first step is to consider what a `distance` function should look like in Python. In other words, what are the inputs (parameters) and what is the output (return value)?

In this case, the two points are the inputs, which we can represent using four parameters. The return value is the distance, which is a floating-point value.

Already we can write an outline of the function that captures our thinking so far:

```
1 def distance(x1, y1, x2, y2):
2     return 0.0
```

Obviously, this version of the function doesn't compute distances; it always returns zero. But it is syntactically correct, and it will run, which means that we can test it before we make it more complicated.

To test the new function, we call it with sample values:

```
>>> distance(1, 2, 4, 6)
0.0
```

We chose these values so that the horizontal distance equals 3 and the vertical distance equals 4; that way, the result is 5 (the hypotenuse of a 3-4-5 triangle). When testing a function, it is useful to know the right answer.

At this point we have confirmed that the function is syntactically correct, and we can start adding lines of code. After each incremental change, we test the function again. If an error occurs at any point, we know where it must be — in the last line we added.

A logical first step in the computation is to find the differences  $x_2 - x_1$  and  $y_2 - y_1$ . We will refer to those values using temporary variables named `dx` and `dy`.

```
1 def distance(x1, y1, x2, y2):
2     dx = x2 - x1
3     dy = y2 - y1
4     return 0.0
```

If we call the function with the arguments shown above, when the flow of execution gets to the return statement, `dx` should be 3 and `dy` should be 4. We can check this by running the function and printing the returned variable.

Next we compute the sum of squares of `dx` and `dy`:

```
1 def distance(x1, y1, x2, y2):
2     dx = x2 - x1
3     dy = y2 - y1
4     dsquared = dx*dx + dy*dy
5     return 0.0
```

Again, we could run the program at this stage and check the value of `dsquared` (which should be 25).

Finally, using the fractional exponent `0.5` to find the square root, we compute and return the result:

```
1 def distance(x1, y1, x2, y2):
2     dx = x2 - x1
3     dy = y2 - y1
4     dsquared = dx*dx + dy*dy
5     result = dsquared**0.5
6     return result
```

If that works correctly, you are done. Otherwise, you might want to inspect the value of `result` before the return statement.

When you start out, you might add only a line or two of code at a time. As you gain more experience, you might find yourself writing and debugging bigger conceptual chunks. Either way, stepping through your code one line at a time and verifying that each step matches your expectations can save you a lot of debugging time. As you improve your programming skills you should find yourself managing bigger and bigger chunks: this is very similar to the way we learned to read letters, syllables, words, phrases, sentences, paragraphs, etc., or the way we learn to chunk music — from individual notes to chords, bars, phrases, and so on.

The key aspects of the process are:

1. Start with a working skeleton program and make small incremental changes. At any point, if there is an error, you will know exactly where it is.
2. Use temporary variables to refer to intermediate values so that you can easily inspect and check them.
3. Once the program is working, relax, sit back, and play around with your options. (There is interesting research that links “playfulness” to better understanding, better learning, more enjoyment, and a more positive mindset about what you can achieve — so spend some time fiddling around!) You might want to consolidate multiple statements into one bigger compound expression, or rename the variables you’ve used, or see if you can make the function shorter. A good guideline is to aim for making code as easy as possible for others to read.

Here is another version of the function. It makes use of a square root function that is in the `math` module (we’ll learn about modules shortly). Which do you prefer? Which looks “closer” to the Pythagorean formula we started out with?

```
1 import math
2
3 def distance(x1, y1, x2, y2):
4     return math.sqrt( (x2-x1)**2 + (y2-y1)**2 )
```



```
>>> distance(1, 2, 4, 6)
5.0
```

## 4.12 Debugging with `print`

A powerful technique for debugging, is to insert extra `print` functions in carefully selected places in your code. Then, by inspecting the output of the program, you can check whether the algorithm is doing what you expect it to. Be clear about the following, however:

- You must have a clear solution to the problem, and must know what should happen before you can debug a program. Work on *solving* the problem on a piece of paper (perhaps using a flowchart to record the steps you take) *before* you concern yourself with writing code. Writing a program doesn't solve the problem — it simply *automates* the manual steps you would take. So first make sure you have a pen-and-paper manual solution that works. Programming then is about making those manual steps happen automatically.
- Do not write **chatterbox** functions. A chatterbox is a fruitful function that, in addition to its primary task, also asks the user for input, or prints output, when it would be more useful if it simply shut up and did its work quietly.

For example, we've seen built-in functions like `range`, `max` and `abs`. None of these would be useful building blocks for other programs if they prompted the user for input, or printed their results while they performed their tasks.

So a good tip is to avoid calling `print` and `input` functions inside fruitful functions, *unless the primary purpose of your function is to perform input and output*. The one exception to this rule might be to temporarily sprinkle some calls to `print` into your code to help debug and understand what is happening when the code runs, but these will then be removed once you get things working.

## 4.13 Composition

As you should expect by now, you can call one function from within another. This ability is called **composition**.

As an example, we'll write a function that takes two points, the center of the circle and a point on the perimeter, and computes the area of the circle.

Assume that the center point is stored in the variables `xc` and `yc`, and the perimeter point is in `xp` and `yp`. The first step is to find the radius of the circle, which is the distance between the two points. Fortunately, we've just written a function, `distance`, that does just that, so now all we have to do is use it:

```
1 radius = distance(xc, yc, xp, yp)
```

The second step is to find the area of a circle with that radius and return it. Again we will use one of our earlier functions:

```
1 result = area(radius)
2 return result
```

Wrapping that up in a function, we get:

```
1 def area_of_circle(xc, yc, xp, yp):
2     radius = distance(xc, yc, xp, yp)
3     result = area(radius)
4     return result
```

The temporary variables `radius` and `result` are useful for development, debugging, and single-stepping through the code to inspect what is happening, but once the program is working, we can make it more concise by composing the function calls:

```
1 def area_of_circle(xc, yc, xp, yp):
2     return area(distance(xc, yc, xp, yp))
```

## 4.14 Boolean functions

Functions can return Boolean values, which is often convenient for hiding complicated tests inside functions. For example:

```
1 def is_divisible(x, y):
2     """ Test if x is exactly divisible by y """
3     if x % y == 0:
4         return True
5     else:
6         return False
```

It is common to give **Boolean functions** names that sound like yes/no questions. `is_divisible` returns either `True` or `False` to indicate whether the `x` is or is not divisible by `y`.

We can make the function more concise by taking advantage of the fact that the condition of the `if` statement is itself a Boolean expression. We can return it directly, avoiding the `if` statement altogether:

```
1 def is_divisible(x, y):
2     return x % y == 0
```

This session shows the new function in action:

```
>>> is_divisible(6, 4)
False
>>> is_divisible(6, 3)
True
```

Boolean functions are often used in conditional statements:

```
1 if is_divisible(x, y):
2     ... # Do something ...
3 else:
4     ... # Do something else ...
```

It might be tempting to write something like:

```
1 if is_divisible(x, y) == True:
```

but the extra comparison is unnecessary.

## 4.15 Programming with style

Readability is very important to programmers, since in practice programs are read and modified far more often than they are written. But, like most rules, we occasionally break them. Most of the code examples in this book will be consistent with the *Python Enhancement Proposal 8* (PEP 8), a style guide developed by the Python community.

We'll have more to say about style as our programs become more complex, but a few pointers will be helpful already:

- use 4 spaces (instead of tabs) for indentation
- limit line length to 78 characters
- when naming identifiers, use `CamelCase` for classes (we'll get to those) and `lowercase_with_underscores` for functions and variables
- place imports at the top of the file
- keep function definitions together below the import statements
- use docstrings to document functions
- use two blank lines to separate function definitions from each other
- keep top level statements, including function calls, together at the bottom of the program

## 4.16 Glossary

**Boolean function** A function that returns a Boolean value. The only possible values of the `bool` type are `False` and `True`.

**chatterbox function** A function which interacts with the user (using `input` or `print`) when it should not. Silent functions that just convert their input arguments into their output results are usually the most useful ones.

**composition (of functions)** Calling one function from within the body of another, or using the return value of one function as an argument to the call of another.

**dead code** Part of a program that can never be executed, often because it appears after a `return` statement.

**fruitful function** A function that yields a return value instead of `None`.

**incremental development** A program development plan intended to simplify debugging by adding and testing only a small amount of code at a time.

**None** A special Python value. One use in Python is that it is returned by functions that do not execute a `return` statement with a return argument.

**return value** The value provided as the result of a function call.

**scaffolding** Code that is used during program development to assist with development and debugging. The unit test code that we added in this chapter are examples of scaffolding.

**temporary variable** A variable used to store an intermediate value in a complex calculation.

## 4.17 Exercises

After completing each exercise, confirm that all the tests pass.

1. The four compass points can be abbreviated by single-letter strings as “N”, “E”, “S”, and “W”. Write a function `turn_clockwise` that takes one of these four compass points as its parameter, and returns the next compass point in the clockwise direction. Here are some tests that should pass:

```
>>>turn_clockwise("N") == "E"
True
>>>turn_clockwise("W") == "N"
True
```

You might ask “What if the argument to the function is some other value?” For all other cases, the function should return the value `None`.

2. Write a function `day_name` that converts an integer number 0 to 6 into the name of a day. Assume day 0 is “Sunday”. Once again, return `None` if the arguments to the function are not valid.
3. Write the inverse function `day_num` which is given a day name, and returns its number.

Once again, if this function is given an invalid argument, it should return `None`.

4. Write a function that helps answer questions like “Today is Wednesday. I leave on holiday in 19 days time. What day will that be?” So the function must take a day name and a `delta` argument — the number of days to add — and should return the resulting day name:

```
day_add("Monday", 4) == "Friday"
day_add("Tuesday", 0) == "Tuesday"
day_add("Tuesday", 14) == "Tuesday"
day_add("Sunday", 100) == "Tuesday"
```

*Hint: use the first two functions written above to help you write this one.*

5. Can your `day_add` function already work with negative deltas? For example, -1 would be yesterday, or -7 would be a week ago:

```
day_add("Sunday", -1) == "Saturday"
day_add("Sunday", -7) == "Sunday"
day_add("Tuesday", -100) == "Sunday"
```

If your function already works, explain why. If it does not work, make it work.

*Hint:* Play with some cases of using the modulus function `%` (introduced at the beginning of the previous chapter). Specifically, explore what happens to `x % 7` when `x` is negative.

6. Write a function `days_in_month` which takes the name of a month, and returns the number of days in the month. Ignore leap years:

```
days_in_month("February") == 28
days_in_month("December") == 31
```

If the function is given invalid arguments, it should return `None`.

7. Write a function `to_secs` that converts hours, minutes and seconds to a total number of seconds. Here are some tests that should pass:

```
to_secs(2, 30, 10) == 9010
to_secs(2, 0, 0) == 7200
to_secs(0, 2, 0) == 120
to_secs(0, 0, 42) == 42
to_secs(0, -10, 10) == -590
```

8. Extend `to_secs` so that it can cope with real values as inputs. It should always return an integer number of seconds (truncated towards zero):

```
to_secs(2.5, 0, 10.71) == 9010
to_secs(2.433, 0, 0) == 8758
```

9. Write three functions that are the “inverses” of `to_secs`:

(a) `hours_in` returns the whole integer number of hours represented by a total number of seconds.

(b) `minutes_in` returns the whole integer number of left over minutes in a total number of seconds, once the hours have been taken out.

(c) `seconds_in` returns the left over seconds represented by a total number of seconds.

You may assume that the total number of seconds passed to these functions is an integer. Here are some test cases:

```
hours_in(9010) == 2
minutes_in(9010) == 30
seconds_in(9010) == 10
```

10. Which of these tests fail? Explain why.

```
3 % 4 == 0
3 % 4 == 3
3 / 4 == 0
3 // 4 == 0
3+4 * 2 == 14
4-2+2 == 0
len("hello, world!") == 13
```

11. Write a compare function that returns 1 if  $a > b$ , 0 if  $a == b$ , and -1 if  $a < b$

```
compare(5, 4) == 1
compare(7, 7) == 0
compare(2, 3) == -1
compare(42, 1) == 1
```

12. Write a function called `hypotenuse` that returns the length of the hypotenuse of a right triangle given the lengths of the two legs as parameters:

```
hypotenuse(3, 4) == 5.0
hypotenuse(12, 5) == 13.0
hypotenuse(24, 7) == 25.0
hypotenuse(9, 12) == 15.0
```

13. Write a function `slope(x1, y1, x2, y2)` that returns the slope of the line through the points  $(x1, y1)$  and  $(x2, y2)$ . Be sure your implementation of `slope` can pass the following tests:

```
slope(5, 3, 4, 2) == 1.0
slope(1, 2, 3, 2) == 0.0
slope(1, 2, 3, 3) == 0.5
slope(2, 4, 1, 2) == 2.0
```

Then use a call to `slope` in a new function named `intercept(x1, y1, x2, y2)` that returns the y-intercept of the line through the points  $(x1, y1)$  and  $(x2, y2)$

```
intercept(1, 6, 3, 12) == 3.0
intercept(6, 1, 1, 6) == 7.0
intercept(4, 6, 12, 8) == 5.0
```

14. Write a function called `is_even(n)` that takes an integer as an argument and returns `True` if the argument is an **even number** and `False` if it is **odd**.

Add your own tests to the test suite.

15. Now write the function `is_odd(n)` that returns `True` when  $n$  is odd and `False` otherwise. Include unit tests for this function too.

Finally, modify it so that it uses a call to `is_even` to determine if its argument is an odd integer, and ensure that its test still pass.

16. Write a function `is_factor(f, n)` that passes these tests:

```
is_factor(3, 12)
not is_factor(5, 12)
is_factor(7, 14)
not is_factor(7, 15)
is_factor(1, 15)
is_factor(15, 15)
not is_factor(25, 15)
```

17. Write `is_multiple` to satisfy these statements using `is_factor` from the previous exercise.

```
is_multiple(12, 3) is_multiple(12, 4) not is_multiple(12, 5) is_multiple(12, 6) not is_multiple(12, 7)
```

18. Write the function `f2c(t)` designed to return the integer value of the nearest degree Celsius for given temperature in Fahrenheit. (*hint*: you may want to make use of the built-in function, `round`. Try printing `round.__doc__` in a Python shell or looking up help for the `round` function, and experimenting with it until you are comfortable with how it works.)

```
f2c(212) == 100      # Boiling point of water
f2c(32) == 0         # Freezing point of water
f2c(-40) == -40      # Wow, what an interesting case!
f2c(36) == 2
f2c(37) == 3
f2c(38) == 3
f2c(39) == 4
```

19. Now do the opposite: write the function `c2f` which converts Celsius to Fahrenheit:

```
c2f(0) == 32
c2f(100) == 212
c2f(-40) == -40
c2f(12) == 54
c2f(18) == 64
c2f(-48) == -54
```

### 4.17.1 Modifiers vs Pure Functions

Functions which take lists as arguments and change them during execution are called **modifiers** and the changes they make are called **side effects**.

A **pure function** does not produce side effects. It communicates with the calling program only through parameters, which it does not modify, and a return value. Let's make a function which doubles the items in a list:

```
1 def double_stuff(values):
2     """ Return a new list which contains
3         doubles of the elements in the list values.
4     """
5     new_list = []
6     for value in values:
7         new_elem = 2 * value
8         new_list.append(new_elem)
9
10    return new_list
```

This version of `double_stuff` does not change its arguments:

```
>>> things = [2, 5, 9]
>>> more_things = double_stuff(things)
>>> things
[2, 5, 9]
>>> more_things
[4, 10, 18]
```

An early rule we saw for assignment said “first evaluate the right hand side, then assign the resulting value to the variable”. So it is quite safe to assign the function result to the same variable that was passed to the function:

```
>>> things = [2, 5, 9]
>>> things = double_stuff(things)
>>> things
[4, 10, 18]
```

If however, we change the definition of `double_stuff` to the following:

```
1 def double_stuff(values):
2     """ Double the elements of values in-place. """
3     for index, value in enumerate(values):
4         values[index] = 2 * value
```

We get upon execution:

```
>>> things = [2, 5, 9]
>>> more_things = double_stuff(things)
>>> things
[4, 10, 18]
>>> more_things
None
```

We see that the original list was modified, while the function doesn’t return anything. This is a good idea when building modifiers.

---

### Which style is better?

In general, we recommend that you always use pure functions, and only use modifiers when you are prepared to stick your head into a lion’s mouth, and have thought about the risks.

---

## 4.17.2 Some Tips, Tricks, and Common Errors

These are small summaries of ideas, tips, and commonly seen errors that might be helpful to those beginning Python.

## 4.18 Functions

Functions help us with our mental chunking: they allow us to group together statements for a high-level purpose, e.g. a function to sort a list of items, a function to make the turtle draw a spiral, or a function to compute the mean and standard deviation of some measurements.

There are two kinds of functions: fruitful, or value-returning functions, which *calculate and return a value*, and we use them because we’re primarily interested in the value they’ll return. Void (non-fruitful) functions are used because they

*perform actions* that we want done — e.g. make a turtle draw a rectangle, or print the first thousand prime numbers. They always return `None` — a special dummy value.

---

**Tip: `None` is not a string**

Values like `None`, `True` and `False` are not strings: they are special values in Python, and are in the list of keywords we gave in chapter 2 (Variables, expressions, and statements). Keywords are special in the language: they are part of the syntax. So we cannot create our own variable or function with a name `True` — we'll get a syntax error. (Built-in functions are not privileged like keywords: we can define our own variable or function called `len`, but we'd be silly to do so!)

---

Along with the fruitful/void families of functions, there are two flavors of the `return` statement in Python: one that returns a useful value, and the other that returns nothing, or `None`. And if we get to the end of any function and we have not explicitly executed any `return` statement, Python automatically returns the value `None`.

---

**Tip: Understand what the function needs to return**

Perhaps nothing — some functions exist purely to perform actions rather than to calculate and return a result. But if the function should return a value, make sure all execution paths do return the value.

---

To make functions more useful, they are given *parameters*. So a function to make a turtle draw a square might have two parameters — one for the turtle that needs to do the drawing, and another for the size of the square. See the first example in Chapter 4 (Functions) — that function can be used with any turtle, and for any size square. So it is much more general than a function that always uses a specific turtle, say `tess` to draw a square of a specific size, say 30.

---

**Tip: Use parameters to generalize functions**

Understand which parts of the function will be hard-coded and unchangeable, and which parts should become parameters so that they can be customized by the caller of the function.

---

---

**Tip: Try to relate Python functions to ideas we already know**

In math, we're familiar with functions like  $f(x) = 3x + 5$ . We already understand that when we call the function  $f(3)$  we make some association between the parameter  $x$  and the argument 3. Try to draw parallels to argument passing in Python.

---

Quiz: Is the function  $f(z) = 3z + 5$  the same as function  $f$  above?

## 4.19 Problems with logic and flow of control

We often want to know if some condition holds for any item in a list, e.g. “does the list have any odd numbers?” This is a common mistake:

```
1 def any_odd(xs): # Buggy version
2     """ Return True if there is an odd number in xs, a list of integers. """
3     for v in xs:
4         if v % 2 == 1:
5             return True
6         else:
7             return False
```



Can we spot two problems here? As soon as we execute a `return`, we'll leave the function. So the logic of saying "If I find an odd number I can return `True`" is fine. However, we cannot return `False` after only looking at one item — we can only return `False` if we've been through all the items, and none of them are odd. So line 6 should not be there, and line 7 has to be outside the loop. To find the second problem above, consider what happens if you call this function with an argument that is an empty list. Here is a corrected version:

```
1 def any_odd(xs):
2     """ Return True if there is an odd number in xs, a list of integers. """
3     for v in xs:
4         if v % 2 == 1:
5             return True
6     return False
```

This "eureka", or "short-circuit" style of returning from a function as soon as we are certain what the outcome will be was first seen in Section 8.10, in the chapter on strings.

It is preferred over this one, which also works correctly:

```
1 def any_odd(xs):
2     """ Return True if there is an odd number in xs, a list of integers. """
3     count = 0
4     for v in xs:
5         if v % 2 == 1:
6             count += 1      # Count the odd numbers
7     if count > 0:
8         return True
9     else:
10        return False
```

The performance disadvantage of this one is that it traverses the whole list, even if it knows the outcome very early on.

---

### Tip: Think about the return conditions of the function

Do I need to look at all elements in all cases? Can I shortcut and take an early exit? Under what conditions? When will I have to examine all the items in the list?

---

The code in lines 7-10 can also be tightened up. The expression `count > 0` evaluates to a Boolean value, either `True` or `False`. The value can be used directly in the `return` statement. So we could cut out that code and simply have the following:

```
1 def any_odd(xs):
2     """ Return True if there is an odd number in xs, a list of integers. """
3     count = 0
4     for v in xs:
5         if v % 2 == 1:
6             count += 1      # Count the odd numbers
7     return count > 0      # Aha! a programmer who understands that Boolean
8                           # expressions are not just used in if statements!
```

Although this code is tighter, it is not as nice as the one that did the short-circuit return as soon as the first odd number was found.

---

### Tip: Generalize your use of Booleans

Mature programmers won't write `if is_prime(n) == True:` when they could say instead `if is_prime(n):`. Think more generally about Boolean values, not just in the context of `if` or `while` statements. Like arithmetic expressions, they have their own set of operators (`and`, `or`, `not`) and values (`True`,

False) and can be assigned to variables, put into lists, etc. A good resource for improving your use of Booleans is [http://en.wikibooks.org/wiki/Non-Programmer%27s\\_Tutorial\\_for\\_Python\\_3/Boolean\\_Expressions](http://en.wikibooks.org/wiki/Non-Programmer%27s_Tutorial_for_Python_3/Boolean_Expressions)

---

Exercise time:

- How would we adapt this to make another function which returns `True` if *all* the numbers are odd? Can you still use a short-circuit style?
- How would we adapt it to return `True` if at least three of the numbers are odd? Short-circuit the traversal when the third odd number is found — don't traverse the whole list unless we have to.

## 4.20 Local variables

Functions are called, or activated, and while they're busy they create their own stack frame which holds local variables. A local variable is one that belongs to the current activation. As soon as the function returns (whether from an explicit return statement or because Python reached the last statement), the stack frame and its local variables are all destroyed. The important consequence of this is that a function cannot use its own variables to remember any kind of state between different activations. It cannot count how many times it has been called, or remember to switch colors between red and blue UNLESS it makes use of variables that are global. Global variables will survive even after our function has exited, so they are the correct way to maintain information between calls.

```
1 sz = 2
2 def h2():
3     """ Draw the next step of a spiral on each call. """
4     global sz
5     tess.turn(42)
6     tess.forward(sz)
7     sz += 1
```

This fragment assumes our turtle is `tess`. Each time we call `h2()` it turns, draws, and increases the global variable `sz`. Python always assumes that an assignment to a variable (as in line 7) means that we want a new local variable, unless we've provided a `global` declaration (on line 4). So leaving out the global declaration means this does not work.

---

### Tip: Local variables do not survive when you exit the function

Use a Python visualizer like the one at [http://netserv.ict.ru.ac.za/python3\\_viz](http://netserv.ict.ru.ac.za/python3_viz) to build a strong understanding of function calls, stack frames, local variables, and function returns.

---

---

### Tip: Assignment in a function creates a local variable

Any assignment to a variable within a function means Python will make a local variable, unless we override with `global`.

---

## 4.21 String handling

There are only four *really* important operations on strings, and we'll be able to do just about anything. There are many more nice-to-have methods (we'll call them sugar coating) that can make life easier, but if we can work with the basic four operations smoothly, we'll have a great grounding.

- `len(str)` finds the length of a string.
- `str[i]` the subscript operation extracts the *i*'th character of the string, as a new string.
- `str[i:j]` the slice operation extracts a substring out of a string.
- `str.find(target)` returns the index where `target` occurs within the string, or -1 if it is not found.

So if we need to know if “snake” occurs as a substring within `s`, we could write

```
1 if s.find("snake") >= 0: ...
2 if "snake" in s: ...           # Also works, nice-to-know sugar coating!
```

It would be wrong to split the string into words unless we were asked whether the *word* “snake” occurred in the string.

Suppose we’re asked to read some lines of data and find function definitions, e.g.: `def some_function_name(x, y):`, and we are further asked to isolate and work with the name of the function. (Let’s say, print it.)

```
1 s = "... "                                # Get the next line from somewhere
2 def_pos = s.find("def ")                  # Look for "def " in the line
3 if def_pos == 0:                          # If it occurs at the left margin
4     op_index = s.find("(")                 # Find the index of the open parenthesis
5     fname = s[4:op_index]                 # Slice out the function name
6     print(fname)                         # ... and work with it.
```

One can extend these ideas:

- What if the function `def` was indented, and didn’t start at column 0? The code would need a bit of adjustment, and we’d probably want to be sure that all the characters in front of the `def_pos` position were spaces. We would not want to do the wrong thing on data like this: `# I def initely like Python!`
- We’ve assumed on line 3 that we will find an open parenthesis. It may need to be checked that we did!
- We have also assumed that there was exactly one space between the keyword `def` and the start of the function name. It will not work nicely for `def f(x)`

As we’ve already mentioned, there are many more “sugar-coated” methods that let us work more easily with strings. There is an `rfind` method, like `find`, that searches from the end of the string backwards. It is useful if we want to find the last occurrence of something. The `lower` and `upper` methods can do case conversion. And the `split` method is great for breaking a string into a list of words, or into a list of lines. We’ve also made extensive use in this book of the `format` method. In fact, if we want to practice reading the Python documentation and learning some new methods on our own, the string methods are an excellent resource.

Exercises:

- Suppose any line of text can contain at most one url that starts with “`http://`” and ends at the next space in the line. Write a fragment of code to extract and print the full url if it is present. (Hint: read the documentation for `find`. It takes some extra arguments, so you can set a starting point from which it will search.)
- Suppose a string contains at most one substring “`< ... >`”. Write a fragment of code to extract and print the portion of the string between the angle brackets.

## 4.22 Looping and lists

Computers are useful because they can repeat computation, accurately and fast. So loops are going to be a central feature of almost all programs you encounter.

---

**Tip: Don’t create unnecessary lists**

Lists are useful if you need to keep data for later computation. But if you don't need lists, it is probably better not to generate them.

---

Here are two functions that both generate ten million random numbers, and return the sum of the numbers. They both work.

```
1  import random
2  joe = random.Random()
3
4  def sum1():
5      """ Build a list of random numbers, then sum them """
6      xs = []
7      for i in range(10000000):
8          num = joe.randrange(1000) # Generate one random number
9          xs.append(num)           # Save it in our list
10
11     tot = sum(xs)
12     return tot
13
14  def sum2():
15      """ Sum the random numbers as we generate them """
16      tot = 0
17      for i in range(10000000):
18          num = joe.randrange(1000)
19          tot += num
20      return tot
21
22  print(sum1())
23  print(sum2())
```

What reasons are there for preferring the second version here? (Hint: open a tool like the Performance Monitor on your computer, and watch the memory usage. How big can you make the list before you get a fatal memory error in `sum1`?)

In a similar way, when working with files, we often have an option to read the whole file contents into a single string, or we can read one line at a time and process each line as we read it. Line-at-a-time is the more traditional and perhaps safer way to do things — you'll be able to work comfortably no matter how large the file is. (And, of course, this mode of processing the files was essential in the old days when computer memories were much smaller.) But you may find whole-file-at-once is sometimes more convenient!



## 5.1 Strings

### 5.1.1 A compound data type

So far we have seen built-in types like `int`, `float`, `bool`, `str` and we've seen lists and pairs. Strings, lists, and pairs are qualitatively different from the others because they are made up of smaller pieces. In the case of strings, they're made up of smaller strings each containing one **character**.

Types that comprise smaller pieces are called **compound data types**. Depending on what we are doing, we may want to treat a compound data type as a single thing, or we may want to access its parts. This ambiguity is useful.

### 5.1.2 Working with strings as single things

We previously saw that each turtle instance has its own attributes and a number of methods that can be applied to the instance. For example, we could set the turtle's color, and we wrote `tess.turn(90)`.

Just like a turtle, a string is also an object. So each string instance has its own attributes and methods.

For example:

```
>>> our_string = "Hello, World!"
>>> all_caps = our_string.upper()
>>> all_caps
'HELLO, WORLD!'
```

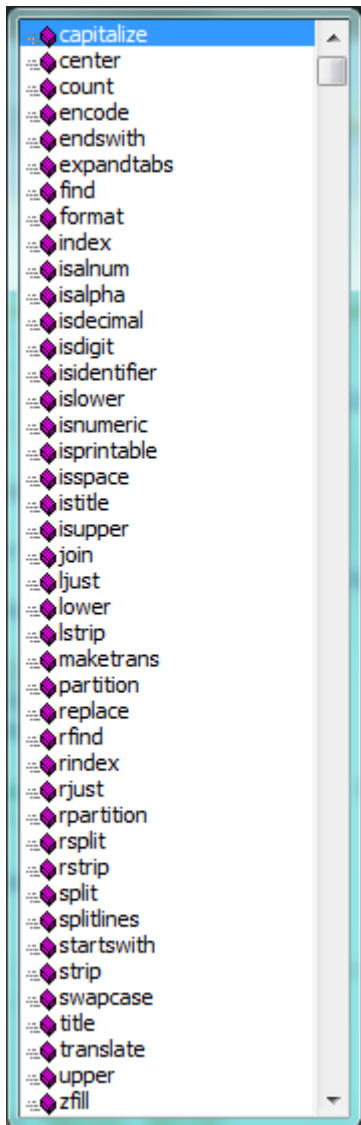
`upper` is a method that can be invoked on any string object to create a new string, in which all the characters are in uppercase. (The original string `our_string` remains unchanged.)

There are also methods such as `lower`, `capitalize`, and `swapcase` that do other interesting stuff.

To learn what methods are available, you can consult the Help documentation, look for string methods, and read the documentation. Or, if you're a bit lazier, simply type the following into an editor like Spyder or PyScripter script:

```
1 our_string = "Hello, World!"
2 new_string = our_string.
```

When you type the period to select one of the methods of `our_string`, your editor might pop up a selection window — typically by pressing `Tab` — showing all the methods (there are around 70 of them — thank goodness we'll only use a few of those!) that could be used on your string.



When you type the name of the method, some further help about its parameter and return type, and its docstring, may be displayed by your scripting environments (for instance, in a Jupyter notebook you can get this information by pressing `Shift+Tab` after a function name).

```
greet = "Hello, World"
xx= greet.swapcase()
print(xx)
```

**\*\* No/Unknown parameters \*\***

`S.swapcase()` -> str

Return a copy of S with uppercase characters converted to lowercase and vice versa.

### 5.1.3 Working with the parts of a string

The **indexing operator** (Python uses square brackets to enclose the index) selects a single character substring from a string:

```
>>> fruit = "banana"
>>> letter = fruit[1]
>>> print(letter)
```

The expression `fruit[1]` selects character number 1 from `fruit`, and creates a new string containing just this one character. The variable `letter` refers to the result. When we display `letter`, we could get a surprise:

```
a
```

Computer scientists always start counting from zero! The letter at subscript position zero of "banana" is b. So at position [1] we have the letter a.

If we want to access the zero-eth letter of a string, we just place 0, or any expression that evaluates to 0, inbetween the brackets:

```
>>> letter = fruit[0]
>>> print(letter)
b
```

The expression in brackets is called an **index**. An index specifies a member of an ordered collection, in this case the collection of characters in the string. The index *indicates* which one you want, hence the name. It can be any integer expression.

We can use `enumerate` to visualize the indices:

```
>>> fruit = "banana"
>>> list(enumerate(fruit))
[(0, 'b'), (1, 'a'), (2, 'n'), (3, 'a'), (4, 'n'), (5, 'a')]
```

Do not worry about `enumerate` at this point, we will see more of it in the chapter on lists.

Note that indexing returns a *string* — Python has no special type for a single character. It is just a string of length 1.

We've also seen lists previously. The same indexing notation works to extract elements from a list:

```
>>> prime_numbers = [2, 3, 5, 7, 11, 13, 17, 19, 23, 29, 31]
>>> prime_numbers[4]
11
>>> friends = ["Joe", "Zoe", "Brad", "Angelina", "Zuki", "Thandi", "Paris"]
>>> friends[3]
'Angelina'
```

### 5.1.4 Length

The `len` function, when applied to a string, returns the number of characters in a string:

```
>>> word = "banana"
>>> len(word)
6
```

To get the last letter of a string, you might be tempted to try something like this:



```
1 size = len(word)
2 last = word[size]           # ERROR!
```

That won't work. It causes the runtime error `IndexError: string index out of range`. The reason is that there is no character at index position 6 in "banana". Because we start counting at zero, the six indexes are numbered 0 to 5. To get the last character, we have to subtract 1 from the length of word:

```
1 size = len(word)
2 last = word[size-1]
```

Alternatively, we can use **negative indices**, which count backward from the end of the string. The expression `word[-1]` yields the last letter, `word[-2]` yields the second to last, and so on.

As you might have guessed, indexing with a negative index also works like this for lists.

### 5.1.5 Traversal and the `for` loop

A lot of computations involve processing a string one character at a time. Often they start at the beginning, select each character in turn, do something to it, and continue until the end. This pattern of processing is called a **traversal**. One way (a very bad way) to encode a traversal is with a `while` statement:

```
1 ix = 0
2 while ix < len(fruit):
3     letter = fruit[ix]
4     print(letter)
5     ix += 1
```

This loop traverses the string and displays each letter on a line by itself. It uses `ix` for the index, which does not make it any clearer. The loop condition is `ix < len(fruit)`, so when `ix` is equal to the length of the string, the condition is false, and the body of the loop is not executed. The last character accessed is the one with the index `len(fruit)-1`, which is the last character in the string. However, this code is a lot longer than it needs to be, and not very clear at all.

But we've previously seen how the `for` loop can easily iterate over the elements in a list and it can do so for strings as well:

```
1 word="Banana"
2 for letter in word:
3     print(letter)
```

Each time through the loop, the next character in the string is assigned to the variable `c`. The loop continues until no characters are left. Here we can see the expressive power the `for` loop gives us compared to the `while` loop when traversing a string.

The following example shows how to use concatenation and a `for` loop to generate an abecedarian series. Abecedarian refers to a series or list in which the elements appear in alphabetical order. For example, in Robert McCloskey's book *Make Way for Ducklings*, the names of the ducklings are Jack, Kack, Lack, Mack, Nack, Ouack, Pack, and Quack. This loop outputs these names in order:

```
1 prefixes = "JKLMNOPQ"
2 suffix = "ack"
3
4 for p in prefixes:
5     print(p + suffix)
```

The output of this program is:

```
Jack
Kack
Lack
Mack
Nack
Oack
Pack
Qack
```

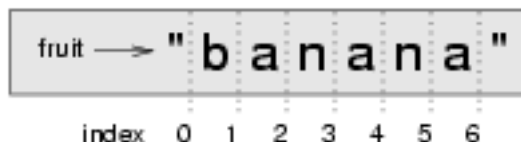
Of course, that's not quite right because Ouack and Quack are misspelled. You'll fix this as an exercise below.

## 5.1.6 Slices

A *substring* of a string is obtained by taking a **slice**. Similarly, we can slice a list to refer to some sublist of the items in the list:

```
>>> phrase = "Pirates of the Caribbean"
>>> print(phrase[0:7])
Pirates
>>> print(phrase[11:14])
the
>>> print(phrase[13:24])
e Caribbean
>>> friends = ["Joe", "Zoe", "Brad", "Angelina", "Zuki", "Thandi", "Paris"]
>>> print(friends[2:4])
['Brad', 'Angelina']
```

The operator `[n:m]` returns the part of the string from the *n*'th character to the *m*'th character, including the first but excluding the last. This behavior makes sense if you imagine the indices pointing *between* the characters, as in the following diagram:



If you imagine this as a piece of paper, the slice operator `[n:m]` copies out the part of the paper between the *n* and *m* positions. Provided *m* and *n* are both within the bounds of the string, your result will be of length (*m*-*n*).

Three tricks are added to this: if you omit the first index (before the colon), the slice starts at the beginning of the string (or list). If you omit the second index, the slice extends to the end of the string (or list). Similarly, if you provide value for *n* that is bigger than the length of the string (or list), the slice will take all the values up to the end. (It won't give an "out of range" error like the normal indexing operation does.) Thus:

```
>>> word = "banana"
>>> word[:3]
'ban'
>>> word[3:]
'ana'
>>> word[3:999]
'ana'
```

What do you think `phrase[:]` means? What about `friends[4:]`? `phrase[-5:-3]`?

### 5.1.7 String comparison

The comparison operators work on strings. To see if two strings are equal:

```
1 if word == "banana":
2     print("Yes, we have no bananas!")
```

Other comparison operations are useful for putting words in *lexicographical* order:

```
1 if word < "banana":
2     print("Your word, " + word + ", comes before banana.")
3 elif word > "banana":
4     print("Your word, " + word + ", comes after banana.")
5 else:
6     print("Yes, we have no bananas!")
```

This is similar to the alphabetical order you would use with a dictionary, except that all the uppercase letters come before all the lowercase letters. As a result:

```
Your word, Zebra, comes before banana.
```

A common way to address this problem is to convert strings to a standard format, such as all lowercase, before performing the comparison. A more difficult problem is making the program realize that zebras are not fruit.

### 5.1.8 Strings are immutable

It is tempting to use the `[]` operator on the left side of an assignment, with the intention of changing a character in a string. For example:

```
1 greeting = "Hello, world!"
2 greeting[0] = 'J'           # ERROR!
3 print(greeting)
```

Instead of producing the output `Jello, world!`, this code produces the runtime error `TypeError: 'str' object does not support item assignment`.

Strings are **immutable**, which means you can't change an existing string. The best you can do is create a new string that is a variation on the original:

```
1 greeting = "Hello, world!"
2 new_greeting = "J" + greeting[1:]
3 print(new_greeting)
```

The solution here is to concatenate a new first letter onto a slice of `greeting`. This operation has no effect on the original string.

### 5.1.9 The `in` and `not in` operators

The `in` operator tests for membership. When both of the arguments to `in` are strings, `in` checks whether the left argument is a substring of the right argument.

```
>>> "p" in "apple"
True
>>> "i" in "apple"
False
```

(continues on next page)

(continued from previous page)

```
>>> "ap" in "apple"
True
>>> "pa" in "apple"
False
```

Note that a string is a substring of itself, and the empty string is a substring of any other string. (Also note that computer scientists like to think about these edge cases quite carefully!)

```
>>> "a" in "a"
True
>>> "apple" in "apple"
True
>>> "" in "a"
True
>>> "" in "apple"
True
```

The `not in` operator returns the logical opposite results of `in`:

```
>>> "x" not in "apple"
True
```

Combining the `in` operator with string concatenation using `+`, we can write a function that removes all the vowels from a string:

```
1 def remove_vowels(phrase):
2     vowels = "aeiou"
3     string_sans_vowels = ""
4     for letter in phrase:
5         if letter.lower() not in vowels:
6             string_sans_vowels += letter
7     return string_sans_vowels
```

Important to note is the `letter.lower()` in line 5, without it, any uppercase vowels would not be removed.

### 5.1.10 A find function

What does the following function do?

```
1 def my_find(haystack, needle):
2     """
3     Find and return the index of needle in haystack.
4     Return -1 if needle does not occur in haystack.
5     """
6     for index, letter in enumerate(haystack):
7         if letter == needle:
8             return index
9     return -1
```

Compare the output of the code above with what Python does itself with the code below:

```
1 haystack = "Bananarama!"
2 print(haystack.find('a'))
3 print(my_find(haystack, 'a'))
```

In a sense, `find` is the opposite of the indexing operator. Instead of taking an index and extracting the corresponding character, it takes a character and finds the index where that character appears. If the character is not found, the function returns `-1`.

This is another example where we see a `return` statement inside a loop. If `letter == needle`, the function returns immediately, breaking out of the loop prematurely.

If the character doesn't appear in the string, then the program exits the loop normally and returns `-1`.

This pattern of computation is sometimes called a **eureka traversal** or **short-circuit evaluation**, because as soon as we find what we are looking for, we can cry "Eureka!", take the short-circuit, and stop looking.

### 5.1.11 Looping and counting

The following program counts the number of times the letter `a` appears in a string, and is another example of the counter pattern introduced in *Counting digits*:

```
1 def count_a(text):
2     count = 0
3     for letter in text:
4         if letter == "a":
5             count += 1
6     return count
7
8 print(count_a("banana") == 3)
```

### 5.1.12 Optional parameters

To find the locations of the second or third occurrence of a character in a string, we can modify the `find` function, adding a third parameter for the starting position in the search string:

```
1 def find2(haystack, needle, start):
2     for index, letter in enumerate(haystack[start:]):
3         if letter == needle:
4             return index + start
5     return -1
6
7
8
9 print(find2("banana", "a", 2) == 3)
```

The call `find2("banana", "a", 2)` now returns 3, the index of the first occurrence of "a" in "banana" starting the search at index 2. What does `find2("banana", "n", 3)` return? If you said, 4, there is a good chance you understand how `find2` works.

Better still, we can combine `find` and `find2` using an **optional parameter**:

```
1 def find(haystack, needle, start=0):
2     for index, letter in enumerate(haystack[start:]):
3         if letter == needle:
4             return index + start
5     return -1
```

When a function has an optional parameter, the caller *may* provide a matching argument. If the third argument is provided to `find`, it gets assigned to `start`. But if the caller leaves the argument out, then `start` is given a default value indicated by the assignment `start=0` in the function definition.

So the call `find("banana", "a", 2)` to this version of `find` behaves just like `find2`, while in the call `find("banana", "a")`, `start` will be set to the **default value** of 0.

Adding another optional parameter to `find` makes it search from a starting position, up to but not including the end position:

```
1 def find(haystack, needle, start=0, end=-1):
2     for index, letter in enumerate(haystack[start:end]):
3         if letter == needle:
4             return index + start
5     return -1
```

The semantics of `start` and `end` in this function are precisely the same as they are in the `range` function.

### 5.1.13 The built-in `find` method

Now that we've done all this work to write a powerful `find` function, we can reveal that strings already have their own built-in `find` method. It can do everything that our code can do, and more! Try all the examples listed above, and check the results!

The built-in `find` method is more general than our version. It can find substrings, not just single characters:

```
>>> "banana".find("nan")
2
>>> "banana".find("na", 3)
4
```

Usually we'd prefer to use the methods that Python provides rather than reinvent our own equivalents. But many of the built-in functions and methods make good teaching exercises, and the underlying techniques you learn are your building blocks to becoming a proficient programmer.

### 5.1.14 The `split` method

One of the most useful methods on strings is the `split` method: it splits a single multi-word string into a list of individual words, removing all the whitespace between them. (Whitespace means any tabs, newlines, or spaces.) This allows us to read input as a single string, and split it into words.

```
>>> phrase = "Well I never did said Alice"
>>> words = phrase.split()
>>> words
['Well', 'I', 'never', 'did', 'said', 'Alice']
```

### 5.1.15 Cleaning up your strings

We'll often work with strings that contain punctuation, or tab and newline characters, especially, as we'll see in a future chapter, when we read our text from files or from the Internet. But if we're writing a program, say, to count word frequencies or check the spelling of each word, we'd prefer to strip off these unwanted characters.

We'll show just one example of how to strip punctuation from a string. Remember that strings are immutable, so we cannot change the string with the punctuation — we need to traverse the original string and create a new string, omitting any punctuation:

```
1 punctuation = "!\"#$%&'()*+,-./:;<=>?@[\\]^_`{|}~"
2
3 def remove_punctuation(phrase):
4     phrase_sans_punct = ""
5     for letter in phrase:
6         if letter not in punctuation:
7             phrase_sans_punct += letter
8     return phrase_sans_punct
```

Setting up that first assignment is messy and error-prone. Fortunately, the Python `string` module already does it for us. So we will make a slight improvement to this program — we'll import the `string` module and use its definition:

```
1 import string
2
3 def remove_punctuation(phrase):
4     phrase_sans_punct = ""
5     for letter in phrase:
6         if letter not in string.punctuation:
7             phrase_sans_punct += letter
8     return phrase_sans_punct
```

Try the examples below: “Well, I never did!”, said Alice. “Are you very, very, sure?”

Composing together this function and the `split` method from the previous section makes a useful combination — we'll clean out the punctuation, and `split` will clean out the newlines and tabs while turning the string into a list of words:

```
1 my_story = """
2 Pythons are constrictors, which means that they will 'squeeze' the life
3 out of their prey. They coil themselves around their prey and with
4 each breath the creature takes the snake will squeeze a little tighter
5 until they stop breathing completely. Once the heart stops the prey
6 is swallowed whole. The entire animal is digested in the snake's
7 stomach except for fur or feathers. What do you think happens to the fur,
8 feathers, beaks, and eggshells? The 'extra stuff' gets passed out as ---
9 you guessed it --- snake POOP! """
10
11 words = remove_punctuation(my_story).split()
12 print(words)
```

The output:

```
['Pythons', 'are', 'constrictors', '...', 'it', 'snake', 'POOP']
```

There are other useful string methods, but this book isn't intended to be a reference manual. On the other hand, the *Python Library Reference* is. Along with a wealth of other documentation, it is available at the [Python website](#).

### 5.1.16 The string format method

The easiest and most powerful way to format a string in Python 3 is to use the `format` method. To see how this works, let's start with a few examples:

```
1 phrase = "His name is {0}!".format("Arthur")
2 print(phrase)
3
4 name = "Alice"
```

(continues on next page)

(continued from previous page)

```
5 age = 10
6 phrase = "I am {1} and I am {0} years old.".format(age, name)
7 print(phrase)
8 phrase = "I am {0} and I am {1} years old.".format(age, name)
9 print(phrase)
10
11 x = 4
12 y = 5
13 phrase = "2**10 = {0} and {1} * {2} = {3:f}".format(2**10, x, y, x * y)
14 print(phrase)
```

Running the script produces:

```
His name is Arthur!
I am Alice and I am 10 years old.
I am 10 and I am Alice years old.
2**10 = 1024 and 4 * 5 = 20.000000
```

The template string contains *place holders*, ... {0} ... {1} ... {2} ... etc. The `format` method substitutes its arguments into the place holders. The numbers in the place holders are indexes that determine which argument gets substituted — make sure you understand line 6 above!

But there's more! Each of the replacement fields can also contain a **format specification** — it is always introduced by the `:` symbol (Line 13 above uses one.) This modifies how the substitutions are made into the template, and can control things like:

- whether the field is aligned to the left `<`, center `^`, or right `>`
- the width allocated to the field within the result string (a number like 10)
- the type of conversion (we'll initially only force conversion to float, `f`, as we did in line 13 of the code above, or perhaps we'll ask integer numbers to be converted to hexadecimal using `x`)
- if the type conversion is a float, you can also specify how many decimal places are wanted (typically, `.2f` is useful for working with currencies to two decimal places.)

Let's do a few simple and common examples that should be enough for most needs. If you need to do anything more esoteric, use `help` and read all the powerful, gory details.

```
1 name1 = "Paris"
2 name2 = "Whitney"
3 name3 = "Hilton"
4
5 print("Pi to three decimal places is {0:.3f}".format(3.1415926))
6 print("123456789 123456789 123456789 123456789 123456789 123456789")
7 print("|||{0:<15}|||{1:^15}|||{2:>15}|||Born in {3}|||")
8     .format(name1,name2,name3,1981))
9 print("The decimal value {0} converts to hex value {0:x}")
10     .format(123456))
```

This script produces the output:

```
Pi to three decimal places is 3.142
123456789 123456789 123456789 123456789 123456789 123456789
|||Paris          |||      Whitney      |||          Hilton|||Born in 1981|||
The decimal value 123456 converts to hex value 1e240
```

You can have multiple placeholders indexing the same argument, or perhaps even have extra arguments that are not referenced at all:



```
1 letter = """
2 Dear {0} {2}.
3 {0}, I have an interesting money-making proposition for you!
4 If you deposit $10 million into my bank account, I can
5 double your money ...
6 """
7
8 print(letter.format("Paris", "Whitney", "Hilton"))
9 print(letter.format("Bill", "Henry", "Gates"))
```

This produces the following:

```
Dear Paris Hilton.
Paris, I have an interesting money-making proposition for you!
If you deposit $10 million into my bank account, I can
double your money ...

Dear Bill Gates.
Bill, I have an interesting money-making proposition for you!
If you deposit $10 million into my bank account I can
double your money ...
```

As you might expect, you'll get an index error if your placeholders refer to arguments that you do not provide:

```
>>> "hello {3}".format("Dave")
Traceback (most recent call last):
  File "<interactive input>", line 1, in <module>
IndexError: tuple index out of range
```

The following example illustrates the real utility of string formatting. First, we'll try to print a table without using string formatting:

```
1 print("i\ti**2\ti**3\ti**5\ti**10\ti**20")
2 for i in range(1, 11):
3     print(i, "\t", i**2, "\t", i**3, "\t", i**5, "\t",
4           i**10, "\t", i**20)
```

This program prints out a table of various powers of the numbers from 1 to 10. (This assumes that the tab width is 8. You might see something even worse than this if you tab width is set to 4.) In its current form it relies on the tab character ( `\t` ) to align the columns of values, but this breaks down when the values in the table get larger than the tab width:

i	i**2	i**3	i**5	i**10	i**20
1	1	1	1	1	1
2	4	8	32	1024	1048576
3	9	27	243	59049	3486784401
4	16	64	1024	1048576	1099511627776
5	25	125	3125	9765625	95367431640625
6	36	216	7776	60466176	3656158440062976
7	49	343	16807	282475249	79792266297612001
8	64	512	32768	1073741824	1152921504606846976
9	81	729	59049	3486784401	12157665459056928801
10	100	1000	100000	10000000000	10000000000000000000

One possible solution would be to change the tab width, but the first column already has more space than it needs. The best solution would be to set the width of each column independently. As you may have guessed by now, string formatting provides a much nicer solution. We can also right-justify each field:

```

1 layout = "{0:>4}{1:>6}{2:>6}{3:>8}{4:>13}{5:>24}"
2
3 print(layout.format("i", "i**2", "i**3", "i**5", "i**10", "i**20"))
4 for i in range(1, 11):
5     print(layout.format(i, i**2, i**3, i**5, i**10, i**20))

```

Running this version produces the following (much more satisfying) output:

i	i**2	i**3	i**5	i**10	i**20
1	1	1	1	1	1
2	4	8	32	1024	1048576
3	9	27	243	59049	3486784401
4	16	64	1024	1048576	1099511627776
5	25	125	3125	9765625	95367431640625
6	36	216	7776	60466176	3656158440062976
7	49	343	16807	282475249	79792266297612001
8	64	512	32768	1073741824	1152921504606846976
9	81	729	59049	3486784401	12157665459056928801
10	100	1000	100000	10000000000	100000000000000000000

## 5.1.17 Summary

This chapter introduced a lot of new ideas. The following summary may prove helpful in remembering what you learned.

**indexing ([ ])** Access a single character in a string using its position (starting from 0). Example: `"This"[2]` evaluates to `"i"`.

**length function (len)** Returns the number of characters in a string. Example: `len("happy")` evaluates to 5.

**for loop traversal (for)** *Traversing* a string means accessing each character in the string, one at a time. For example, the following for loop:

```

for ch in "Example":
    ...

```

executes the body of the loop 7 times with different values of `ch` each time.

**slicing ([ : ])** A *slice* is a substring of a string. Example: `'bananas and cream'[3:6]` evaluates to `ana` (so does `'bananas and cream'[1:4]`).

**string comparison (>, <, >=, <=, ==, !=)** The six common comparison operators work with strings, evaluating according to *lexicographical* order. Examples: `"apple" < "banana"` evaluates to `True`. `"Zeta" < "Apricot"` evaluates to `False`. `"Zebra" <= "aardvark"` evaluates to `True` because all upper case letters precede lower case letters.

**in and not in operator (in, not in)** The `in` operator tests for membership. In the case of strings, it tests whether one string is contained inside another string. Examples: `"heck" in "I'll be checking for you."` evaluates to `True`. `"cheese" in "I'll be checking for you."` evaluates to `False`.

## 5.1.18 Glossary

**compound data type** A data type in which the values are made up of components, or elements, that are themselves values.

**default value** The value given to an optional parameter if no argument for it is provided in the function call.

**docstring** A string constant on the first line of a function or module definition (and as we will see later, in class and method definitions as well). Docstrings provide a convenient way to associate documentation with code. Docstrings are also used by programming tools to provide interactive help.

**dot notation** Use of the **dot operator**, `.`, to access methods and attributes of an object.

**immutable data value** A data value which cannot be modified. Assignments to elements or slices (sub-parts) of immutable values cause a runtime error.

**index** A variable or value used to select a member of an ordered collection, such as a character from a string, or an element from a list.

**mutable data value** A data value which can be modified. The types of all mutable values are compound types. Lists and dictionaries are mutable; strings and tuples are not.

**optional parameter** A parameter written in a function header with an assignment to a default value which it will receive if no corresponding argument is given for it in the function call.

**short-circuit evaluation** A style of programming that shortcuts extra work as soon as the outcome is known with certainty. In this chapter our `find` function returned as soon as it found what it was looking for; it didn't traverse all the rest of the items in the string.

**slice** A part of a string (substring) specified by a range of indices. More generally, a subsequence of any sequence type in Python can be created using the slice operator (`sequence[start:stop]`).

**traverse** To iterate through the elements of a collection, performing a similar operation on each.

**whitespace** Any of the characters that move the cursor without printing visible characters. The constant `string.whitespace` contains all the white-space characters.

### 5.1.19 Exercises

1. What is the result of each of the following:

```
>>> "Python"[1]
>>> "Strings are sequences of characters."[5]
>>> len("wonderful")
>>> "Mystery"[:4]
>>> "p" in "Pineapple"
>>> "apple" in "Pineapple"
>>> "pear" not in "Pineapple"
>>> "apple" > "pineapple"
>>> "pineapple" < "Peach"
```

2. Modify:

```
1 prefixes = "JKLMNOPQ"
2 suffix = "ack"
3
4 for letter in prefixes:
5     print(letter + suffix)
```

so that Ouack and Quack are spelled correctly.

3. Encapsulate

```
1 word = "banana"
2 count = 0
3 for letter in word:
```

(continues on next page)

(continued from previous page)

```

4     if letter == "a":
5         count += 1
6     print(count)

```

in a function named `count_letters`, and generalize it so that it accepts the string and the letter as arguments. Make the function return the number of characters, rather than print the answer. The caller should do the printing.

- Now rewrite the `count_letters` function so that instead of traversing the string, it repeatedly calls the `find` method, with the optional third parameter to locate new occurrences of the letter being counted.
- Assign to a variable in your program a triple-quoted string that contains your favourite paragraph of text — perhaps a poem, a speech, instructions to bake a cake, some inspirational verses, etc.

Write a function which removes all punctuation from the string, breaks the string into a list of words, and counts the number of words in your text that contain the letter “e”. Your program should print an analysis of the text like this:

```
Your text contains 243 words, of which 109 (44.8%) contain an "e".
```

- Print a neat looking multiplication table like this:

	1	2	3	4	5	6	7	8	9	10	11	12
1:	1	2	3	4	5	6	7	8	9	10	11	12
2:	2	4	6	8	10	12	14	16	18	20	22	24
3:	3	6	9	12	15	18	21	24	27	30	33	36
4:	4	8	12	16	20	24	28	32	36	40	44	48
5:	5	10	15	20	25	30	35	40	45	50	55	60
6:	6	12	18	24	30	36	42	48	54	60	66	72
7:	7	14	21	28	35	42	49	56	63	70	77	84
8:	8	16	24	32	40	48	56	64	72	80	88	96
9:	9	18	27	36	45	54	63	72	81	90	99	108
10:	10	20	30	40	50	60	70	80	90	100	110	120
11:	11	22	33	44	55	66	77	88	99	110	121	132
12:	12	24	36	48	60	72	84	96	108	120	132	144

- Write a function that reverses its string argument, and satisfies these tests:

```

1 reverse("happy") == "yppah"
2 reverse("Python") == "nohtyP"
3 reverse("") == ""
4 reverse("a") == "a"

```

- Write a function that mirrors its argument:

```

1 mirror("good") == "gooddoog"
2 mirror("Python") == "PythonnohtyP"
3 mirror("") == ""
4 mirror("a") == "aa"

```

- Write a function that removes all occurrences of a given letter from a string:

```

1 remove_letter("a", "apple") == "pple"
2 remove_letter("a", "banana") == "bnn"
3 remove_letter("z", "banana") == "banana"
4 remove_letter("i", "Mississippi") == "Msssspp"
5 remove_letter("b", "") == ""
6 remove_letter("b", "c") == "c"

```

10. Write a function that recognizes palindromes. (Hint: use your `reverse` function to make this easy!):

```
1 is_palindrome("abba")
2 not is_palindrome("abab")
3 is_palindrome("tenet")
4 not is_palindrome("banana")
5 is_palindrome("straw warts")
6 is_palindrome("a")
7 # is_palindrome("")      # Is an empty string a palindrome?
```

11. Write a function that counts how many times a substring occurs in a string:

```
1 count("is", "Mississippi") == 2
2 count("an", "banana") == 2
3 count("ana", "banana") == 2
4 count("nana", "banana") == 1
5 count("nanan", "banana") == 0
6 count("aaa", "aaaaaa") == 4
```

12. Write a function that removes the first occurrence of a string from another string:

```
1 remove("an", "banana") == "bana"
2 remove("cyc", "bicycle") == "bile"
3 remove("iss", "Mississippi") == "Missippi"
4 remove("eggs", "bicycle") == "bicycle"
```

13. Write a function that removes all occurrences of a string from another string:

```
1 remove_all("an", "banana") == "ba"
2 remove_all("cyc", "bicycle") == "bile"
3 remove_all("iss", "Mississippi") == "Mippi"
4 remove_all("eggs", "bicycle") == "bicycle"
```

There are only four *really* important operations on strings, and we'll be able to do just about anything. There are many more nice-to-have methods (we'll call them sugar coating) that can make life easier, but if we can work with the basic four operations smoothly, we'll have a great grounding.

- `len(str)` finds the length of a string.
- `str[i]` the subscript operation extracts the *i*'th character of the string, as a new string.
- `str[i:j]` the slice operation extracts a substring out of a string.
- `str.find(target)` returns the index where `target` occurs within the string, or -1 if it is not found.

So if we need to know if "snake" occurs as a substring within `s`, we could write

```
1 if s.find("snake") >= 0: ...
2 if "snake" in s: ...      # Also works, nice-to-know sugar coating!
```

It would be wrong to split the string into words unless we were asked whether the *word* "snake" occurred in the string.

Suppose we're asked to read some lines of data and find function definitions, e.g.: `def some_function_name(x, y):`, and we are further asked to isolate and work with the name of the function. (Let's say, print it.)

```
1 s = "... "                # Get the next line from somewhere
2 def_pos = s.find("def ")   # Look for "def " in the line
3 if def_pos == 0:           # If it occurs at the left margin
4     op_index = s.find("(")  # Find the index of the open parenthesis
```

(continues on next page)

(continued from previous page)

```
5     fnname = s[4:op_index]           # Slice out the function name
6     print(fnname)                   # ... and work with it.
```

One can extend these ideas:

- What if the function def was indented, and didn't start at column 0? The code would need a bit of adjustment, and we'd probably want to be sure that all the characters in front of the `def_pos` position were spaces. We would not want to do the wrong thing on data like this: `# I def initely like Python!`
- We've assumed on line 3 that we will find an open parenthesis. It may need to be checked that we did!
- We have also assumed that there was exactly one space between the keyword `def` and the start of the function name. It will not work nicely for `def f(x)`

As we've already mentioned, there are many more “sugar-coated” methods that let us work more easily with strings. There is an `rfind` method, like `find`, that searches from the end of the string backwards. It is useful if we want to find the last occurrence of something. The `lower` and `upper` methods can do case conversion. And the `split` method is great for breaking a string into a list of words, or into a list of lines. We've also made extensive use in this book of the `format` method. In fact, if we want to practice reading the Python documentation and learning some new methods on our own, the string methods are an excellent resource.

Exercises:

- Suppose any line of text can contain at most one url that starts with “`http://`” and ends at the next space in the line. Write a fragment of code to extract and print the full url if it is present. (Hint: read the documentation for `find`. It takes some extra arguments, so you can set a starting point from which it will search.)
- Suppose a string contains at most one substring “`< ... >`”. Write a fragment of code to extract and print the portion of the string between the angle brackets.

## 5.2 Tuples

### 5.2.1 Tuples are used for grouping data

We saw earlier that we could group together pairs of values by surrounding with parentheses. Recall this example:

```
>>> year_born = ("Paris Hilton", 1981)
```

This is an example of a **data structure** — a mechanism for grouping and organizing data to make it easier to use.

The pair is an example of a **tuple**. Generalizing this, a tuple can be used to group any number of items into a single compound value. Syntactically, a tuple is a comma-separated sequence of values. Although it is not necessary, it is conventional to enclose tuples in parentheses:

```
>>> julia = ("Julia", "Roberts", 1967, "Duplicity", 2009, "Actress",
↳ "Atlanta, Georgia")
```

The other thing that could be said somewhere around here, is that the parentheses are there to disambiguate. For example, if we have a tuple nested within another tuple and the parentheses weren't there, how would we tell where the nested tuple begins/ends? Also: the creation of an empty tuple is done like this: `empty_tuple=()`

Tuples are useful for representing what other languages often call *records* (or structs) — some related information that belongs together, like your student record. There is no description of what each of these fields means, but we can guess. A tuple lets us “chunk” together related information and use it as a single thing.

Tuples support the same sequence operations as strings. The index operator selects an element from a tuple.

```
>>> julia[2]
1967
```

But if we try to use item assignment to modify one of the elements of the tuple, we get an error:

```
>>> julia[0] = "X"
TypeError: 'tuple' object does not support item assignment
```

So like strings, tuples are immutable. Once Python has created a tuple in memory, it cannot be changed.

Of course, even if we can’t modify the elements of a tuple, we can always make the `julia` variable reference a new tuple holding different information. To construct the new tuple, it is convenient that we can slice parts of the old tuple and join up the bits to make the new tuple. So if `julia` has a new recent film, we could change her variable to reference a new tuple that used some information from the old one:

```
>>> julia = julia[:3] + ("Eat Pray Love", 2010) + julia[5:]
>>> julia
('Julia', 'Roberts', 1967, 'Eat Pray Love', 2010, 'Actress', 'Atlanta, ↵
↵Georgia')
```

To create a tuple with a single element (but you’re probably not likely to do that too often), we have to include the final comma, because without the final comma, Python treats the (5) below as an integer in parentheses:

```
>>> tup = (5,)
>>> type(tup)
<class 'tuple'>
>>> x = (5)
>>> type(x)
<class 'int'>
```

### 5.2.2 Tuple assignment

Python has a very powerful **tuple assignment** feature that allows a tuple of variables on the left of an assignment to be assigned values from a tuple on the right of the assignment. (We already saw this used for pairs, but it generalizes.)

```
(name, surname, year_born, movie, year_movie, profession, birthplace) = julia
```

This does the equivalent of seven assignment statements, all on one easy line. One requirement is that the number of variables on the left must match the number of elements in the tuple.

One way to think of tuple assignment is as tuple packing/unpacking.

In tuple packing, the values on the left are ‘packed’ together in a tuple:

```
>>> bob = ("Bob", 19, "CS")      # tuple packing
```

In tuple unpacking, the values in a tuple on the right are ‘unpacked’ into the variables/names on the right:

```
>>> bob = ("Bob", 19, "CS")
>>> (name, age, studies) = bob    # tuple unpacking
>>> name
```

(continues on next page)

(continued from previous page)

```
'Bob'  
>>> age  
19  
>>> studies  
'CS'
```

Once in a while, it is useful to swap the values of two variables. With conventional assignment statements, we have to use a temporary variable. For example, to swap `a` and `b`:

```
1 temp = a  
2 a = b  
3 b = temp
```

Tuple assignment solves this problem neatly:

```
1 (a, b) = (b, a)
```

The left side is a tuple of variables; the right side is a tuple of values. Each value is assigned to its respective variable. All the expressions on the right side are evaluated before any of the assignments. This feature makes tuple assignment quite versatile.

Naturally, the number of variables on the left and the number of values on the right have to be the same:

```
>>> (one, two, three, four) = (1, 2, 3)  
ValueError: need more than 3 values to unpack
```

### 5.2.3 Tuples as return values

Functions can always only return a single value, but by making that value a tuple, we can effectively group together as many values as we like, and return them together. This is very useful — we often want to know some batsman's highest and lowest score, or we want to find the mean and the standard deviation, or we want to know the year, the month, and the day, or if we're doing some ecological modelling we may want to know the number of rabbits and the number of wolves on an island at a given time.

For example, we could write a function that returns both the area and the circumference of a circle of radius `r`:

```
1 def circle_stats(r):  
2     """ Return (circumference, area) of a circle of radius r """  
3     circumference = 2 * math.pi * r  
4     area = math.pi * r * r  
5     return (circumference, area)
```

### 5.2.4 Composability of Data Structures

We saw in an earlier chapter that we could make a list of pairs, and we had an example where one of the items in the tuple was itself a list:

```
students = [  
    ("John", ["CompSci", "Physics"]),  
    ("Vusi", ["Maths", "CompSci", "Stats"]),  
    ("Jess", ["CompSci", "Accounting", "Economics", "Management"]),  
    ("Sarah", ["InfSys", "Accounting", "Economics", "CommLaw"]),  
    ("Zuki", ["Sociology", "Economics", "Law", "Stats", "Music"])]
```



Tuples items can themselves be other tuples. For example, we could improve the information about our movie stars to hold the full date of birth rather than just the year, and we could have a list of some of her movies and dates that they were made, and so on:

```
julia_more_info = ( ("Julia", "Roberts"), (8, "October", 1967),
                    "Actress", ("Atlanta", "Georgia"),
                    [ ("Duplicity", 2009),
                      ("Notting Hill", 1999),
                      ("Pretty Woman", 1990),
                      ("Erin Brockovich", 2000),
                      ("Eat Pray Love", 2010),
                      ("Mona Lisa Smile", 2003),
                      ("Oceans Twelve", 2004) ])
```

Notice in this case that the tuple has just five elements — but each of those in turn can be another tuple, a list, a string, or any other kind of Python value. This property is known as being **heterogeneous**, meaning that it can be composed of elements of different types.

### 5.2.5 Glossary

**data structure** An organization of data for the purpose of making it easier to use.

**immutable data value** A data value which cannot be modified. Assignments to elements or slices (sub-parts) of immutable values cause a runtime error.

**mutable data value** A data value which can be modified. The types of all mutable values are compound types. Lists and dictionaries are mutable; strings and tuples are not.

**tuple** An immutable data value that contains related elements. Tuples are used to group together related data, such as a person's name, their age, and their gender.

**tuple assignment** An assignment to all of the elements in a tuple using a single assignment statement. Tuple assignment occurs *simultaneously* rather than in sequence, making it useful for swapping values.

### 5.2.6 Exercises

1. We've said nothing in this chapter about whether you can pass tuples as arguments to a function. Construct a small Python example to test whether this is possible, and write up your findings.
2. Is a pair a generalization of a tuple, or is a tuple a generalization of a pair?
3. Is a pair a kind of tuple, or is a tuple a kind of pair?

## 5.3 Lists

A **list** is an ordered collection of values. The values that make up a list are called its **elements**, or its **items**. We will use the term *element* or *item* to mean the same thing. Lists are similar to strings, which are ordered collections of characters, except that the elements of a list can be of any type. Lists and strings — and other collections that maintain the order of their items — are called **sequences**.

### 5.3.1 List values

There are several ways to create a new list; the simplest is to enclose the elements in square brackets ([ and ]):

```
1 numbers = [10, 20, 30, 40]
2 words = ["spam", "bungee", "swallow"]
```

The first example is a list of four integers. The second is a list of three strings. The elements of a list don't have to be the same type. The following list contains a string, a float, an integer, and (amazingly) another list:

```
1 stuffs = ["hello", 2.0, 5, [10, 20]]
```

A list within another list is said to be **nested**.

Finally, a list with no elements is called an empty list, and is denoted `[]`.

We have already seen that we can assign list values to variables or pass lists as parameters to functions:

```
1 >>> vocabulary = ["apple", "cheese", "dog"]
2 >>> numbers = [17, 123]
3 >>> an_empty_list = []
4 >>> print(vocabulary, numbers, an_empty_list)
5 ["apple", "cheese", "dog"] [17, 123] []
```

## 5.3.2 Accessing elements

The syntax for accessing the elements of a list is the same as the syntax for accessing the characters of a string — the index operator: `[]` (not to be confused with an empty list). The expression inside the brackets specifies the index. Remember that the indices start at 0:

```
>>> numbers[0]
17
```

Any expression evaluating to an integer can be used as an index:

```
>>> numbers[9-8]
123
>>> numbers[1.0]
Traceback (most recent call last):
  File "<interactive input>", line 1, in <module>
TypeError: list indices must be integers, not float
```

If you try to access or assign to an element that does not exist, you get a runtime error:

```
>>> numbers[2]
Traceback (most recent call last):
  File "<interactive input>", line 1, in <module>
IndexError: list index out of range
```

It is common (but wrong!) to use a loop variable as a list index.

```
1 horsemen = ["war", "famine", "pestilence", "death"]
2
3 for i in [0, 1, 2, 3]:
4     print(horsemen[i])
```

Each time through the loop, the variable `i` is used as an index into the list, printing the `i`'th element. This pattern of computation is called a **list traversal**.

The above sample doesn't need or use the index `i` for anything besides getting the items from the list, so this more direct version — where the `for` loop gets the items — is much more clear!

```
1 horsemen = ["war", "famine", "pestilence", "death"]
2
3 for h in horsemen:
4     print(h)
```

### 5.3.3 List length

The function `len` returns the length of a list, which is equal to the number of its elements. If you are going to use an integer index to access the list, it is a good idea to use this value as the upper bound of a loop instead of a constant. That way, if the size of the list changes, you won't have to go through the program changing all the loops; they will work correctly for any size list:

```
1 horsemen = ["war", "famine", "pestilence", "death"]
2
3 for i in range(len(horsemen)):
4     print(horsemen[i])
```

The last time the body of the loop is executed, `i` is `len(horsemen) - 1`, which is the index of the last element. (But the version without the index looks even better now! The version above is not the right way to do things!)

```
1 horsemen = ["war", "famine", "pestilence", "death"]
2
3 for horseman in horsemen:
4     print(horseman)
```

Although a list can contain another list, the nested list still counts as a single element in its parent list. The length of this list is 4:

```
>>> len(["car makers", 1, ["Ford", "Toyota", "BMW"], [1, 2, 3]])
4
```

### 5.3.4 List membership

`in` and `not in` are Boolean operators that test membership in a sequence. We used them previously with strings, but they also work with lists and other sequences:

```
>>> horsemen = ["war", "famine", "pestilence", "death"]
>>> "pestilence" in horsemen
True
>>> "debauchery" in horsemen
False
>>> "debauchery" not in horsemen
True
```

Using this produces a more elegant version of the nested loop program we previously used to count the number of students doing Computer Science in the section *Nested Loops for Nested Data*:

```
1 students = [
2     ("John", ["CompSci", "Physics"]),
3     ("Vusi", ["Maths", "CompSci", "Stats"]),
4     ("Jess", ["CompSci", "Accounting", "Economics", "Management"]),
5     ("Sarah", ["InfSys", "Accounting", "Economics", "CommLaw"]),
6     ("Zuki", ["Sociology", "Economics", "Law", "Stats", "Music"])]
```

(continues on next page)

(continued from previous page)

```
7
8 # Count how many students are taking CompSci
9 counter = 0
10 for name, subjects in students:
11     if "CompSci" in subjects:
12         counter += 1
13
14 print("The number of students taking CompSci is", counter)
```

### 5.3.5 List operations

The + operator concatenates lists:

```
>>> first_list = [1, 2, 3]
>>> second_list = [4, 5, 6]
>>> both_lists = first_list + second_list
>>> both_lists
[1, 2, 3, 4, 5, 6]
```

Similarly, the \* operator repeats a list a given number of times:

```
>>> [0] * 4
[0, 0, 0, 0]
>>> [1, 2, 3] * 3
[1, 2, 3, 1, 2, 3, 1, 2, 3]
```

The first example repeats [0] four times. The second example repeats the list [1, 2, 3] three times.

### 5.3.6 List slices

The slice operations we saw previously with strings let us work with sublists:

```
>>> a_list = ["a", "b", "c", "d", "e", "f"]
>>> a_list[1:3]
['b', 'c']
>>> a_list[:4]
['a', 'b', 'c', 'd']
>>> a_list[3:]
['d', 'e', 'f']
>>> a_list[:]
['a', 'b', 'c', 'd', 'e', 'f']
```

### 5.3.7 Lists are mutable

Unlike strings, lists are **mutable**, which means we can change their elements. Using the index operator on the left side of an assignment, we can update one of the elements:

```
>>> fruit = ["banana", "apple", "quince"]
>>> fruit[0] = "pear"
>>> fruit[2] = "orange"
>>> fruit
['pear', 'apple', 'orange']
```

The bracket operator applied to a list can appear anywhere in an expression. When it appears on the left side of an assignment, it changes one of the elements in the list, so the first element of `fruit` has been changed from "banana" to "pear", and the last from "quince" to "orange". An assignment to an element of a list is called **item assignment**. Item assignment does not work for strings:

```
>>> my_string = "TEST"
>>> my_string[2] = "X"
Traceback (most recent call last):
  File "<interactive input>", line 1, in <module>
TypeError: 'str' object does not support item assignment
```

but it does for lists:

```
>>> my_list = ["T", "E", "S", "T"]
>>> my_list[2] = "X"
>>> my_list
['T', 'E', 'X', 'T']
```

With the slice operator we can update a whole sublist at once:

```
>>> a_list = ["a", "b", "c", "d", "e", "f"]
>>> a_list[1:3] = ["x", "y"]
>>> a_list
['a', 'x', 'y', 'd', 'e', 'f']
```

We can also remove elements from a list by assigning an empty list to them:

```
>>> a_list = ["a", "b", "c", "d", "e", "f"]
>>> a_list[1:3] = []
>>> a_list
['a', 'd', 'e', 'f']
```

And we can add elements to a list by squeezing them into an empty slice at the desired location:

```
>>> a_list = ["a", "d", "f"]
>>> a_list[1:1] = ["b", "c"]
>>> a_list
['a', 'b', 'c', 'd', 'f']
>>> a_list[4:4] = ["e"]
>>> a_list
['a', 'b', 'c', 'd', 'e', 'f']
```

### 5.3.8 List deletion

Using slices to delete list elements can be error-prone. Python provides an alternative that is more readable. The `del` statement removes an element from a list:

```
>>> a = ["one", "two", "three"]
>>> del a[1]
>>> a
['one', 'three']
```

As you might expect, `del` causes a runtime error if the index is out of range.

You can also use `del` with a slice to delete a sublist:

```
>>> a_list = ["a", "b", "c", "d", "e", "f"]
>>> del a_list[1:5]
>>> a_list
['a', 'f']
```

As usual, the sublist selected by slice contains all the elements up to, but not including, the second index.

### 5.3.9 Objects and references

After we execute these assignment statements

```
1 a = "banana"
2 b = "banana"
```

we know that `a` and `b` will refer to a string object with the letters `"banana"`. But we don't know yet whether they point to the *same* string object.

There are two possible ways the Python interpreter could arrange its memory:



In one case, `a` and `b` refer to two different objects that have the same value. In the second case, they refer to the same object.

We can test whether two names refer to the same object using the `is` operator:

```
>>> a is b
True
```

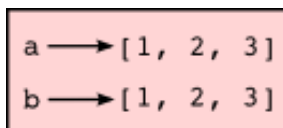
This tells us that both `a` and `b` refer to the same object, and that it is the second of the two state snapshots that accurately describes the relationship.

Since strings are *immutable*, Python optimizes resources by making two names that refer to the same string value refer to the same object.

This is not the case with lists:

```
>>> a = [1, 2, 3]
>>> b = [1, 2, 3]
>>> a == b
True
>>> a is b
False
```

The state snapshot here looks like this:



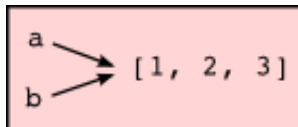
`a` and `b` have the same value but do not refer to the same object.

### 5.3.10 Aliasing

Since variables refer to objects, if we assign one variable to another, both variables refer to the same object:

```
>>> a = [1, 2, 3]
>>> b = a
>>> a is b
True
```

In this case, the state snapshot looks like this:



Because the same list has two different names, *a* and *b*, we say that it is **aliased**. Changes made with one alias affect the other:

```
>>> b[0] = 5
>>> a
[5, 2, 3]
```

Although this behavior can be useful, it is sometimes unexpected or undesirable. In general, it is safer to avoid aliasing when you are working with mutable objects (i.e. lists at this point in our textbook, but we'll meet more mutable objects as we cover classes and objects, dictionaries and sets). Of course, for immutable objects (i.e. strings, tuples), there's no problem — it is just not possible to change something and get a surprise when you access an alias name. That's why Python is free to alias strings (and any other immutable kinds of data) when it sees an opportunity to economize.

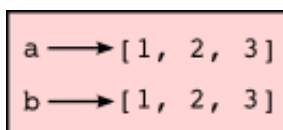
### 5.3.11 Cloning lists

If we want to modify a list and also keep a copy of the original, we need to be able to make a copy of the list itself, not just the reference. This process is sometimes called **cloning**, to avoid the ambiguity of the word copy.

The easiest way to clone a list is to use the slice operator:

```
>>> a = [1, 2, 3]
>>> b = a[:]
>>> b
[1, 2, 3]
```

Taking any slice of *a* creates a new list. In this case the slice happens to consist of the whole list. So now the relationship is like this:



Now we are free to make changes to *b* without worrying that we'll inadvertently be changing *a*:

```
>>> b[0] = 5
>>> a
[1, 2, 3]
```

### 5.3.12 Lists and for loops

The `for` loop also works with lists, as we've already seen. The generalized syntax of a `for` loop is:

```
for <VARIABLE> in <LIST>:
    <BODY>
```

So, as we've seen

```
1 friends = ["Joe", "Zoe", "Brad", "Angelina", "Zuki", "Thandi", "Paris"]
2 for friend in friends:
3     print(friend)
```

It almost reads like English: For (every) friend in (the list of) friends, print (the name of the) friend.

Any list expression can be used in a `for` loop:

```
1 for number in range(20):
2     if number % 3 == 0:
3         print(number)
4
5 for fruit in ["banana", "apple", "quince"]:
6     print("I like to eat " + fruit + "s!")
```

The first example prints all the multiples of 3 between 0 and 19. The second example expresses enthusiasm for various fruits.

Since lists are mutable, we often want to traverse a list, changing each of its elements. The following squares all the numbers in the list `xs`:

```
1 xs = [1, 2, 3, 4, 5]
2
3 for i in range(len(xs)):
4     xs[i] = xs[i]**2
```

Take a moment to think about `range(len(xs))` until you understand how it works.

In this example we are interested in both the *value* of an item, (we want to square that value), and its *index* (so that we can assign the new value to that position). This pattern is common enough that Python provides a nicer way to implement it:

```
1 xs = [1, 2, 3, 4, 5]
2
3 for (i, val) in enumerate(xs):
4     xs[i] = val**2
```

`enumerate` generates pairs of both (index, value) during the list traversal. Try this next example to see more clearly how `enumerate` works:

```
1 for (i, v) in enumerate(["banana", "apple", "pear", "lemon"]):
2     print(i, v)
```

```
0 banana
1 apple
2 pear
3 lemon
```



### 5.3.13 List parameters

Passing a list as an argument actually passes a reference to the list, not a copy or clone of the list. So parameter passing creates an alias for you: the caller has one variable referencing the list, and the called function has an alias, but there is only one underlying list object. For example, the function below takes a list as an argument and multiplies each element in the list by 2:

```
1 def double_stuff(stuff_list):
2     """ Overwrite each element in a_list with double its value. """
3     for (index, stuff) in enumerate(stuff_list):
4         stuff_list[index] = 2 * stuff
```

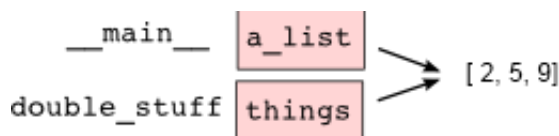
If we add the following onto our script:

```
1 things = [2, 5, 9]
2 double_stuff(things)
3 print(things)
```

When we run it we'll get:

```
[4, 10, 18]
```

In the function above, the parameter `stuff_list` and the variable `things` are aliases for the same object. So before any changes to the elements in the list, the state snapshot looks like this:



Since the list object is shared by two frames, we drew it between them.

If a function modifies the items of a list parameter, the caller sees the change.

### 5.3.14 List methods

The dot operator can also be used to access built-in methods of list objects. We'll start with the most useful method for adding something onto the end of an existing list:

```
>>> mylist = []
>>> mylist.append(5)
>>> mylist.append(27)
>>> mylist.append(3)
>>> mylist.append(12)
>>> mylist
[5, 27, 3, 12]
```

`append` is a list method which adds the argument passed to it to the end of the list. We'll use it heavily when we're creating new lists. Continuing with this example, we show several other list methods:

```
>>> mylist.insert(1, 12)  # Insert 12 at pos 1, shift other items up
>>> mylist
[5, 12, 27, 3, 12]
>>> mylist.count(12)      # How many times is 12 in mylist?
2
>>> mylist.extend([5, 9, 5, 11])  # Put whole list onto end of mylist
```

(continues on next page)

(continued from previous page)

```
>>> mylist
[5, 12, 27, 3, 12, 5, 9, 5, 11]
>>> mylist.index(9)           # Find index of first 9 in mylist
6
>>> mylist.reverse()
>>> mylist
[11, 5, 9, 5, 12, 3, 27, 12, 5]
>>> mylist.sort()
>>> mylist
[3, 5, 5, 5, 9, 11, 12, 12, 27]
>>> mylist.remove(12)         # Remove the first 12 in the list
>>> mylist
[3, 5, 5, 5, 9, 11, 12, 27]
```

Experiment and play with the list methods shown here, and read their documentation until you feel confident that you understand how they work.

### 5.3.15 Pure functions and modifiers

As seen before, there is a difference between a pure function and one with side-effects. The difference is shown below as lists have some special gotcha's. Functions which take lists as arguments and change them during execution are called **modifiers** and the changes they make are called **side effects**.

A **pure function** does not produce side effects. It communicates with the calling program only through parameters, which it does not modify, and a return value. Here is `double_stuff` written as a pure function:

```
1 def double_stuff(a_list):
2     """ Return a new list which contains
3         doubles of the elements in a_list.
4     """
5     new_list = []
6     for value in a_list:
7         new_elem = 2 * value
8         new_list.append(new_elem)
9
10    return new_list
```

This version of `double_stuff` does not change its arguments:

```
>>> things = [2, 5, 9]
>>> more_things = double_stuff(things)
>>> things
[2, 5, 9]
>>> more_things
[4, 10, 18]
```

An early rule we saw for assignment said “first evaluate the right hand side, then assign the resulting value to the variable”. So it is quite safe to assign the function result to the same variable that was passed to the function:

```
>>> things = [2, 5, 9]
>>> things = double_stuff(things)
>>> things
[4, 10, 18]
```

### 5.3.16 Functions that produce lists

The pure version of `double_stuff` above made use of an important **pattern** for your toolbox. Whenever you need to write a function that creates and returns a list, the pattern is usually:

```
1 initialize a result variable to be an empty list
2 loop
3     create a new element
4     append it to result
5 return the result
```

Let us show another use of this pattern. Assume you already have a function `is_prime(x)` that can test if `x` is prime. Write a function to return a list of all prime numbers less than `n`:

```
1 def primes_less_than(n):
2     """ Return a list of all prime numbers less than n. """
3     result = []
4     for i in range(2, n):
5         if is_prime(i):
6             result.append(i)
7     return result
```

### 5.3.17 Strings and lists

Two of the most useful methods on strings involve conversion to and from lists of substrings. The `split` method (which we've already seen) breaks a string into a list of words. By default, any number of whitespace characters is considered a word boundary:

```
>>> song = "The rain in Spain..."
>>> words = song.split()
>>> words
['The', 'rain', 'in', 'Spain...']
```

An optional argument called a **delimiter** can be used to specify which string to use as the boundary marker between substrings. The following example uses the string `ai` as the delimiter:

```
>>> song.split("ai")
['The r', 'n in Sp', 'n...']
```

Notice that the delimiter doesn't appear in the result.

The inverse of the `split` method is `join`. You choose a desired **separator** string, (often called the *glue*) and join the list with the glue between each of the elements:

```
>>> glue = ";"
>>> phrase = glue.join(words)
>>> phrase
'The;rain;in;Spain...'
```

The list that you glue together (`words` in this example) is not modified. Also, as these next examples show, you can use empty glue or multi-character strings as glue:

```
>>> " --- ".join(words)
'The --- rain --- in --- Spain...'
>>> "".join(words)
'TheraininSpain...'
```

### 5.3.18 list and range

Python has a built-in type conversion function called `list` that tries to turn whatever you give it into a list.

```
>>> letters = list("Crunchy Frog")
>>> letters
['C', 'r', 'u', 'n', 'c', 'h', 'y', ' ', 'F', 'r', 'o', 'g']
>>> "".join(letters)
'Crunchy Frog'
```

One particular feature of `range` is that it doesn't instantly compute all its values: it “puts off” the computation, and does it on demand, or “lazily”. We'll say that it gives a **promise** to produce the values when they are needed. This is very convenient if your computation short-circuits a search and returns early, as in this case:

```
1 def f(n):
2     """ Find the first positive integer between 101 and less
3         than n that is divisible by 21
4     """
5     for i in range(101, n):
6         if (i % 21 == 0):
7             return i
8
9
10 print(f(110) == 105)
11 print(f(1000000000) == 105)
```

In the second test, if `range` were to eagerly go about building a list with all those elements, you would soon exhaust your computer's available memory and crash the program. But it is cleverer than that! This computation works just fine, because the `range` object is just a promise to produce the elements if and when they are needed. Once the condition in the `if` becomes true, no further elements are generated, and the function returns. (Note: Before Python 3, `range` was not lazy. If you use an earlier versions of Python, YMMV!)

---

#### YMMV: Your Mileage May Vary

The acronym YMMV stands for *your mileage may vary*. American car advertisements often quoted fuel consumption figures for cars, e.g. that they would get 28 miles per gallon. But this always had to be accompanied by legal small-print warning the reader that they might not get the same. The term YMMV is now used idiomatically to mean “your results may differ”, e.g. *The battery life on this phone is 3 days, but YMMV*.

---

You'll sometimes find the lazy `range` wrapped in a call to `list`. This forces Python to turn the lazy promise into an actual list:

```
>>> range(10)                # Create a lazy promise
range(0, 10)
>>> list(range(10))          # Call in the promise, to produce a list.
[0, 1, 2, 3, 4, 5, 6, 7, 8, 9]
```

### 5.3.19 Looping and lists

Computers are useful because they can repeat computation, accurately and fast. So loops are going to be a central feature of almost all programs you encounter.

---

**Tip: Don't create unnecessary lists**

Lists are useful if you need to keep data for later computation. But if you don't need lists, it is probably better not to generate them.

---

Here are two functions that both generate ten million random numbers, and return the sum of the numbers. They both work.

```
1 import random
2 joe = random.Random()
3
4 def sum1():
5     """ Build a list of random numbers, then sum them """
6     xs = []
7     for i in range(10000000):
8         num = joe.randrange(1000) # Generate one random number
9         xs.append(num)           # Save it in our list
10
11     tot = sum(xs)
12     return tot
13
14 def sum2():
15     """ Sum the random numbers as we generate them """
16     tot = 0
17     for i in range(10000000):
18         num = joe.randrange(1000)
19         tot += num
20     return tot
21
22 print(sum1())
23 print(sum2())
```

What reasons are there for preferring the second version here? (Hint: open a tool like the Performance Monitor on your computer, and watch the memory usage. How big can you make the list before you get a fatal memory error in `sum1`?)

In a similar way, when working with files, we often have an option to read the whole file contents into a single string, or we can read one line at a time and process each line as we read it. Line-at-a-time is the more traditional and perhaps safer way to do things — you'll be able to work comfortably no matter how large the file is. (And, of course, this mode of processing the files was essential in the old days when computer memories were much smaller.) But you may find whole-file-at-once is sometimes more convenient!

### 5.3.20 Nested lists

A nested list is a list that appears as an element in another list. In this list, the element with index 3 is a nested list:

```
>>> nested = ["hello", 2.0, 5, [10, 20]]
```

If we output the element at index 3, we get:

```
>>> print(nested[3])
[10, 20]
```

To extract an element from the nested list, we can proceed in two steps:

```
>>> elem = nested[3]
>>> elem[0]
10
```

Or we can combine them:

```
>>> nested[3][1]
20
```

Bracket operators evaluate from left to right, so this expression gets the 3'th element of `nested` and extracts the 1'th element from it.

### 5.3.21 Matrices

Nested lists are often used to represent matrices. For example, the matrix:

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$$

might be represented as:

```
>>> mx = [[1, 2, 3], [4, 5, 6], [7, 8, 9]]
```

`mx` is a list with three elements, where each element is a row of the matrix. We can select an entire row from the matrix in the usual way:

```
>>> mx[1]
[4, 5, 6]
```

Or we can extract a single element from the matrix using the double-index form:

```
>>> mx[1][2]
6
```

The first index selects the row, and the second index selects the column. Although this way of representing matrices is common, it is not the only possibility. A small variation is to use a list of columns instead of a list of rows. Later we will see a more radical alternative using a dictionary.

### 5.3.22 Glossary

**aliases** Multiple variables that contain references to the same object.

**clone** To create a new object that has the same value as an existing object. Copying a reference to an object creates an alias but doesn't clone the object.

**delimiter** A character or string used to indicate where a string should be split.

**element** One of the values in a list (or other sequence). The bracket operator selects elements of a list. Also called *item*.

**immutable data value** A data value which cannot be modified. Assignments to elements or slices (sub-parts) of immutable values cause a runtime error.

**index** An integer value that indicates the position of an item in a list. Indexes start from 0.

**item** See *element*.

**list** A collection of values, each in a fixed position within the list. Like other types `str`, `int`, `float`, etc. there is also a `list` type-converter function that tries to turn whatever argument you give it into a list.

**list traversal** The sequential accessing of each element in a list.

**modifier** A function which changes its arguments inside the function body. Only mutable types can be changed by modifiers.

**mutable data value** A data value which can be modified. The types of all mutable values are compound types. Lists and dictionaries are mutable; strings and tuples are not.

**nested list** A list that is an element of another list.

**object** A thing to which a variable can refer.

**pattern** A sequence of statements, or a style of coding something that has general applicability in a number of different situations. Part of becoming a mature Computer Scientist is to learn and establish the patterns and algorithms that form your toolkit. Patterns often correspond to your “mental chunking”.

**promise** An object that promises to do some work or deliver some values if they’re eventually needed, but it lazily puts off doing the work immediately. Calling `range` produces a promise.

**pure function** A function which has no side effects. Pure functions only make changes to the calling program through their return values.

**sequence** Any of the data types that consist of an ordered collection of elements, with each element identified by an index.

**side effect** A change in the state of a program made by calling a function. Side effects can only be produced by modifiers.

**step size** The interval between successive elements of a linear sequence. The third (and optional argument) to the `range` function is called the step size. If not specified, it defaults to 1.

### 5.3.23 Exercises

1. What is the Python interpreter’s response to the following?

```
>>> list(range(10, 0, -2))
```

The three arguments to the `range` function are *start*, *stop*, and *step*, respectively. In this example, *start* is greater than *stop*. What happens if *start* < *stop* and *step* < 0? Write a rule for the relationships among *start*, *stop*, and *step*.

2. Consider this fragment of code:

```
1 import turtle
2
3 tess = turtle.Turtle()
4 alex = tess
5 alex.color("hotpink")
```

Does this fragment create one or two turtle instances? Does setting the color of `alex` also change the color of `tess`? Explain in detail.

3. Draw a state snapshot for `a` and `b` before and after the third line of the following Python code is executed:

```
1 a = [1, 2, 3]
2 b = a[:]
3 b[0] = 5
```

4. What will be the output of the following program?

```

1 this = ["I", "am", "not", "a", "crook"]
2 that = ["I", "am", "not", "a", "crook"]
3 print("Test 1: {}".format(this is that))
4 that = this
5 print("Test 2: {}".format(this is that))

```

Provide a *detailed* explanation of the results.

5. Lists can be used to represent mathematical *vectors*. In this exercise and several that follow you will write functions to perform standard operations on vectors. Create a script named `vectors.py` and write Python code to pass the tests in each case.

Write a function `add_vectors(vector1, vector2)` that takes two lists of numbers of the same length, and returns a new list containing the sums of the corresponding elements of each:

```

1 add_vectors([1, 1], [1, 1]) == [2, 2]
2 add_vectors([1, 2], [1, 4]) == [2, 6]
3 add_vectors([1, 2, 1], [1, 4, 3]) == [2, 6, 4]

```

6. Write a function `scalar_mult(scalar, vector)` that takes a number, `scalar`, and a list, `vector` and returns the *scalar multiple* of `vector` by `scalar`.

```

1 scalar_mult(5, [1, 2]) == [5, 10]
2 scalar_mult(3, [1, 0, -1]) == [3, 0, -3]
3 scalar_mult(7, [3, 0, 5, 11, 2]) == [21, 0, 35, 77, 14]

```

7. Write a function `dot_product(vec1, vec2)` that takes two lists of numbers of the same length, and returns the sum of the products of the corresponding elements of each (the *dot product*).

```

1 dot_product([1, 1], [1, 1]) == 2
2 dot_product([1, 2], [1, 4]) == 9
3 dot_product([1, 2, 1], [1, 4, 3]) == 12

```

8. *Extra challenge for the mathematically inclined:* Write a function `cross_product(vec1, vec2)` that takes two lists of numbers of length 3 and returns their *cross product*. You should write your own tests.
9. Describe the relationship between `" ".join(song.split())` and `song` in the fragment of code below. Are they the same for all strings assigned to `song`? When would they be different?

```

1 song = "The rain in Spain..."

```

10. Write a function `replace(s, old, new)` that replaces all occurrences of `old` with `new` in a string `s`:

```

1 replace("Mississippi", "i", "I") == "MIssIssIppI"
2
3 song = "I love spom! Spom is my favorite food. Spom, spom, yum!"
4 replace(song, "om", "am") ==
5     "I love spam! Spam is my favorite food. Spam, spam, yum!"
6
7 replace(song, "o", "a") ==
8     "I lave spam! Spam is my favarite faad. Spam, spam, yum!"

```

*Hint:* use the `split` and `join` methods.



## 5.4 Dictionaries

All of the compound data types we have studied in detail so far — strings, lists, and tuples — are sequence types, which use integers as indices to access the values they contain within them.

**Dictionaries** are yet another kind of compound type. They are Python's built-in **mapping type**. They map **keys**, which can be any immutable type, to values, which can be any type (heterogeneous), just like the elements of a list or tuple. In other languages, they are called associative arrays since they associate a key with a value.

As an example, we will create a dictionary to translate English words into Spanish. For this dictionary, the keys are strings.

One way to create a dictionary is to start with the empty dictionary and add **key:value pairs**. The empty dictionary is denoted `{}`:

```
>>> english_spanish = {}
>>> english_spanish["one"] = "uno"
>>> english_spanish["two"] = "dos"
```

The first assignment creates a dictionary named `english_spanish`; the other assignments add new key:value pairs to the dictionary. We can print the current value of the dictionary in the usual way:

```
>>> print(english_spanish)
{'two': 'dos', 'one': 'uno'}
```

The key:value pairs of the dictionary are separated by commas. Each pair contains a key and a value separated by a colon.

---

### Hashing

The order of the pairs may not be what was expected. Python uses complex algorithms, designed for very fast access, to determine where the key:value pairs are stored in a dictionary. For our purposes we can think of this ordering as unpredictable.

You also might wonder why we use dictionaries at all when the same concept of mapping a key to a value could be implemented using a list of tuples:

```
>>> {"apples": 430, "bananas": 312, "oranges": 525, "pears": 217}
{'pears': 217, 'apples': 430, 'oranges': 525, 'bananas': 312}
>>> [('apples', 430), ('bananas', 312), ('oranges', 525), ('pears', 217)]
[('apples', 430), ('bananas', 312), ('oranges', 525), ('pears', 217)]
```

The reason is dictionaries are very fast, implemented using a technique called hashing, which allows us to access a value very quickly. By contrast, the list of tuples implementation is slow. If we wanted to find a value associated with a key, we would have to iterate over every tuple, checking the 0th element. What if the key wasn't even in the list? We would have to get to the end of it to find out.

---

Another way to create a dictionary is to provide a list of key:value pairs using the same syntax as the previous output:

```
>>> english_spanish = {"one": "uno", "two": "dos", "three": "tres"}
```

It doesn't matter what order we write the pairs. The values in a dictionary are accessed with keys, not with indices, so there is no need to care about ordering.

Here is how we use a key to look up the corresponding value:

```
>>> print(english_spanish["two"])
'dos'
```

The key "two" yields the value "dos".

Lists, tuples, and strings have been called *sequences*, because their items occur in order. The dictionary is the first compound type that we've seen that is not a sequence, so we can't index or slice a dictionary.

### 5.4.1 Dictionary operations

The `del` statement removes a key:value pair from a dictionary. For example, the following dictionary contains the names of various fruits and the number of each fruit in stock:

```
>>> inventory = {"apples": 430, "bananas": 312, "oranges": 525, "pears": 217}
>>> print(inventory)
{'pears': 217, 'apples': 430, 'oranges': 525, 'bananas': 312}
```

If someone buys all of the bananas, we can remove the entry from the dictionary:

```
>>> del inventory["bananas"]
>>> print(inventory)
{'apples': 430, 'oranges': 525, 'pears': 217}
```

If we then try to see how many bananas we have, we get an error (because, yes, we have no bananas). (Try this!)

Or if we're expecting more bananas soon, we might just change the value associated with bananas:

```
>>> inventory["bananas"] = 0
>>> print(inventory)
{'pears': 217, 'apples': 430, 'oranges': 525, 'bananas': 0}
```

A new shipment of bananas arriving could be handled like this:

```
>>> inventory["bananas"] += 200
>>> print(inventory)
{'pears': 0, 'apples': 430, 'oranges': 525, 'bananas': 512}
```

The `len` function also works on dictionaries; it returns the number of key:value pairs:

```
>>> len(inventory)
4
```

### 5.4.2 Dictionary methods

Dictionaries have a number of useful built-in methods.

The `keys` method returns what Python 3 calls a **view** of its underlying keys. A view object has some similarities to the `range` object we saw earlier — it is a lazy promise, to deliver its elements when they're needed by the rest of the program. We can iterate over the view, or turn the view into a list like this:

```
1 for key in english_spanish.keys(): # The order of the k's is not defined
2     print("Got key", key, "which maps to value", english_spanish[key])
3
4 keys = list(english_spanish.keys())
5 print(keys)
```

This produces this output:

```
Got key three which maps to value tres
Got key two which maps to value dos
Got key one which maps to value uno
['three', 'two', 'one']
```

It is so common to iterate over the keys in a dictionary that we can omit the `keys` method call in the `for` loop — iterating over a dictionary implicitly iterates over its keys:

```
1 for key in english_spanish:
2     print("Got key", key)
```

The `values` method is similar; it returns a view object which can be turned into a list:

```
>>> list(english_spanish.values())
['tres', 'dos', 'uno']
```

The `items` method also returns a view, which promises a list of tuples — one tuple for each key:value pair:

```
>>> list(english_spanish.items())
[('three', 'tres'), ('two', 'dos'), ('one', 'uno')]
```

Tuples are often useful for getting both the key and the value at the same time while we are looping:

```
1 for (key,value) in english_spanish.items():
2     print("Got",key,"that maps to",value)
```

This produces:

```
Got three that maps to tres
Got two that maps to dos
Got one that maps to uno
```

The `in` and `not in` operators can test if a key is in the dictionary:

```
>>> "one" in english_spanish
True
>>> "six" in english_spanish
False
>>> "tres" in english_spanish    # Note that 'in' tests keys, not values.
False
```

This method can be very useful, since looking up a non-existent key in a dictionary causes a runtime error:

```
>>> english_spanish["dog"]
Traceback (most recent call last):
...
KeyError: 'dog'
```

## 5.4.3 Aliasing and copying

As in the case of lists, because dictionaries are mutable, we need to be aware of aliasing. Whenever two variables refer to the same object, changes to one affect the other.

If we want to modify a dictionary and keep a copy of the original, use the `copy` method. For example, `opposites` is a dictionary that contains pairs of opposites:

```
>>> opposites = {"up": "down", "right": "wrong", "yes": "no"}
>>> alias = opposites
>>> copy = opposites.copy()    # Shallow copy
```

`alias` and `opposites` refer to the same object; `copy` refers to a fresh copy of the same dictionary. If we modify `alias`, `opposites` is also changed:

```
>>> alias["right"] = "left"
>>> opposites["right"]
'left'
```

If we modify `copy`, `opposites` is unchanged:

```
>>> copy["right"] = "privilege"
>>> opposites["right"]
'left'
```

## 5.4.4 Counting letters

In the exercises in Chapter 8 (Strings) we wrote a function that counted the number of occurrences of a letter in a string. A more general version of this problem is to form a frequency table of the letters in the string, that is, how many times each letter appears.

Such a frequency table might be useful for compressing a text file. Because different letters appear with different frequencies, we can compress a file by using shorter codes for common letters and longer codes for letters that appear less frequently.

Dictionaries provide an elegant way to generate a frequency table:

```
>>> letter_counts = {}
>>> for letter in "Mississippi":
...     letter_counts[letter] = letter_counts.get(letter, 0) + 1
...
>>> letter_counts
{'M': 1, 's': 4, 'p': 2, 'i': 4}
```

We start with an empty dictionary. For each letter in the string, we find the current count (possibly zero) and increment it. At the end, the dictionary contains pairs of letters and their frequencies.

It might be more appealing to display the frequency table in alphabetical order. We can do that with the `items` and `sort` methods (more precisely, `sort` orders lexicographically):

```
>>> letter_items = list(letter_counts.items())
>>> letter_items.sort()
>>> print(letter_items)
[('M', 1), ('i', 4), ('p', 2), ('s', 4)]
```

Notice in the first line we had to call the type conversion function `list`. That turns the promise we get from `items` into a list, a step that is needed before we can use the list's `sort` method.

## 5.4.5 Glossary

**call graph** A graph consisting of nodes which represent function frames (or invocations), and directed edges (lines with arrows) showing which frames gave rise to other frames.

**dictionary** A collection of key:value pairs that maps from keys to values. The keys can be any immutable value, and the associated value can be of any type.

**immutable data value** A data value which cannot be modified. Assignments to elements or slices (sub-parts) of immutable values cause a runtime error.

**key** A data item that is *mapped to* a value in a dictionary. Keys are used to look up values in a dictionary. Each key must be unique across the dictionary.

**key:value pair** One of the pairs of items in a dictionary. Values are looked up in a dictionary by key.

**mapping type** A mapping type is a data type comprised of a collection of keys and associated values. Python's only built-in mapping type is the dictionary. Dictionaries implement the [associative array](#) abstract data type.

**memo** Temporary storage of precomputed values to avoid duplicating the same computation.

**mutable data value** A data value which can be modified. The types of all mutable values are compound types. Lists and dictionaries are mutable; strings and tuples are not.

## 5.4.6 Exercises

1. Write a program that reads a string and returns a table of the letters of the alphabet in alphabetical order which occur in the string together with the number of times each letter occurs. Case should be ignored. A sample output of the program when the user enters the data "ThiS is String with Upper and lower case Letters", would look like this:

```
a  2
c  1
d  1
e  5
g  1
h  2
i  4
l  2
n  2
o  1
p  2
r  4
s  5
t  5
u  1
w  2
```

2. Give the Python interpreter's response to each of the following from a continuous interpreter session:

(a) 

```
>>> dictionary = {"apples": 15, "bananas": 35, "grapes": 12}
>>> dictionary["bananas"]
```

(b) 

```
>>> dictionary["oranges"] = 20
>>> len(dictionary)
```

(c) 

```
>>> "grapes" in dictionary
```

(d) 

```
>>> dictionary["pears"]
```

(e) 

```
>>> dictionary.get("pears", 0)
```

```
(f) >>> fruits = list(dictionary.keys())
>>> fruits.sort()
>>> print(fruits)
```

```
(g) >>> del dictionary["apples"]
>>> "apples" in dictionary
```

Be sure you understand why you get each result. Then apply what you have learned to fill in the body of the function below:

```
1 def add_fruit(inventory, fruit, quantity=0):
2     return None
3
4 # Make these tests work...
5 new_inventory = {}
6 add_fruit(new_inventory, "strawberries", 10)
7 print("strawberries" in new_inventory)
8 print(new_inventory["strawberries"] == 10)
9 add_fruit(new_inventory, "strawberries", 25)
10 print(new_inventory["strawberries"] == 35)
```

3. Write a program called `alice_words.py` that creates a text file named `alice_words.txt` containing an alphabetical listing of all the words, and the number of times each occurs, in the text version of *Alice's Adventures in Wonderland*. (You can obtain a free plain text version of the book, along with many others, from <http://www.gutenberg.org>.) The first 10 lines of your output file should look something like this:

Word	Count
a	631
a-piece	1
abide	1
able	1
about	94
above	3
absence	1
absurd	2

How many times does the word `alice` occur in the book?

4. What is the longest word in *Alice in Wonderland*? How many characters does it have?



## CHAPTER 6

---

### Numpy

---

The standard Python data types are not very suited for mathematical operations. For example, suppose we have the list `a = [2, 3, 8]`. If we multiply this list by an integer, we get:

```
>>> a = [2, 3, 8]
>>> 2 * a
[2, 3, 8, 2, 3, 8]
```

And float's are not even allowed:

```
>>> a = [2, 3, 8]
>>> 2 * a
>>> 2.1 * a
TypeError: can't multiply sequence by non-int of type 'float'
```

In order to solve this using Python lists, we would have to do something like:

```
values = [2, 3, 8]
result = []
for x in values:
    result.append(2.1 * x)
```

This is not very elegant, is it? This is because Python `list`'s are not designed as mathematical objects. Rather, they are purely a collection of items. In order to get a type of list which behaves like a mathematical array or matrix, we use Numpy.

```
>>> import numpy as np
>>> a = np.array([2, 3, 8])
>>> 2.1 * a
array([ 4.2,  6.3, 16.8])
```

As we can see, this worked the way we expected it to. We note a couple of things: - We abbreviated numpy to np, this is conventional. - `np.array` takes a Python list as argument. - The list `[2, 3, 8]` contains `int`'s, yet the result contains `float`'s. This means numpy changed the data type automatically for us.

Now let's take it a step further and see what happens when we multiply together array's.



```
>>> import numpy as np
>>> a = np.array([2, 3, 8])
>>> a * a
array([ 4,  9, 64])
>>> a**2
array([ 4,  9, 64])
```

This has nicely squared the array element-wise.

---

**Note:** Those in the know might be a bit surprised by this. After all, if  $a$  is a vector, shouldn't  $a**2$  be the dot product of the two vectors,  $\vec{a} \cdot \vec{a}$ ? Well, `numpy` arrays are not vectors in the algebraic sense. Arithmetic operations between arrays are performed element-wise, not on the arrays as a whole.

To tell `numpy` we want the dot product we simply use the `np.dot` function:

```
>>> a = np.array([2, 3, 8])
>>> np.dot(a,a)
77
```

Furthermore, if you pass 2D arrays to `np.dot` it will behave like matrix multiplication. Several other similar NumPy algebraic functions are available (like `np.cross`, `np.outer`, etc.)

**Bottom line:** when you want to treat `numpy` array operations as vector or matrix operations, make use of the specialized functions to this end.

---

## 6.1 Shape

One of the most important properties an array is its shape. We have already seen 1 dimensional (1D) arrays, but arrays can have any dimensions you like. Images for example, consist of a 2D array of pixels. But in color images every pixel is an RGB tuple: the intensity in red, green and blue. Every pixel itself is therefore an array as well. This makes a color image 3D overall.

To get the shape of an array, we use `shape`:

```
>>> import numpy as np
>>> a = np.array([2, 3, 8])
>>> a.shape
(3,)
```

Something slightly more interesting:

```
>>> b = np.array([
    [2, 3, 8],
    [4, 5, 6],
])
>>> b.shape
(2, 3)
```

## 6.2 Slicing

Just like with lists, we might want to select certain values from an array. For 1D arrays it works just like for normal python lists:

```
>>> a = np.array([2, 3, 8])
>>> a[2]
8
>>> a[1:]
np.array([3, 8])
```

However, when dealing with higher dimensional arrays something else happens:

```
>>> b = np.array([
    [2, 3, 8],
    [4, 5, 6],
])
>>> b[1]
array([4, 5, 6])
>>> b[1][2]
6
```

We see that using `b[1]` returns the 1th row along the first dimension, which is still an array. After that, we can select individual items from that. This can be abbreviated to:

```
>>> b[1, 2]
6
```

But what if I wanted the 1th column instead of the first row? Then we use `:` to select all items along the first dimension, and then a `1`:

```
>>> b[:, 1]
array([3, 5])
```

By comparing with the definition of `b`, we see that this is the column we were looking for.

---

**Note:** Instead of first, I write 1th on purpose to signify the existence of a 0th element. Remember that in Python, as in any self-respecting programming language, we start counting at zero.

---

Find out more about advanced slicing at the [Numpy indexing documentation](#) page.

## 6.3 Masking

This is perhaps the single most powerful feature of Numpy. Suppose we have an array, and we want to throw away all values above a certain cutoff:

```
>>> a = np.array([230, 10, 284, 39, 76])
>>> cutoff = 200
>>> a > cutoff
np.array([True, False, True, False, False])
```

Simply using the larger than operator lets us know in which cases the test was positive. Now we set all the values above 200 to zero:

```
>>> a = np.array([230, 10, 284, 39, 76])
>>> cutoff = 200
>>> a[a > cutoff] = 0
>>> a
np.array([0, 10, 0, 39, 76])
```

The crucial line is `a[a > cutoff] = 0`. This selects all the points in the array where the test was positive and assigns 0 to that position. Without knowing this trick we would have had to loop over the array:

```
>>> a = np.array([230, 10, 284, 39, 76])
>>> cutoff = 200
>>> new_a = []
>>> for x in a:
>>>     if x > cutoff:
>>>         new_a.append(0)
>>>     else:
>>>         new_a.append(x)
>>> a = np.array(new_a)
```

Looks rather silly now, doesn't it? When working with images this becomes even more obvious, because there we might have to loop over three dimensions before we can use the if/else. Can you imagine the mess?

## 6.4 Broadcasting

Another powerful feature of Numpy is broadcasting. Broadcasting takes place when you perform operations between arrays of different shapes. For instance

```
>>> a = np.array([
    [0, 1],
    [2, 3],
    [4, 5],
    ])
>>> b = np.array([10, 100])
>>> a * b
array([[ 0, 100],
       [ 20, 300],
       [ 40, 500]])
```

The shapes of `a` and `b` don't match. In order to proceed, Numpy will stretch `b` into a second dimension, as if it were stacked three times upon itself. The operation then takes place element-wise.

One of the rules of broadcasting is that only dimensions of size 1 can be stretched (if an array only has one dimension, all other dimensions are considered for broadcasting purposes to have size 1). In the example above `b` is 1D, and has shape (2,). For broadcasting with `a`, which has two dimensions, Numpy adds another dimension of size 1 to `b`. `b` now has shape (1, 2). This new dimension can now be stretched three times so that `b`'s shape matches `a`'s shape of (3, 2).

The other rule is that dimensions are compared from the last to the first. Any dimensions that do not match must be stretched to become equally sized. However, according to the previous rule, only dimensions of size 1 can stretch. This means that some shapes cannot broadcast and Numpy will give you an error:

```
>>> c = np.array([
    [0, 1, 2],
    [3, 4, 5],
    ])
>>> b = np.array([10, 100])
>>> c * b
ValueError: operands could not be broadcast together with shapes (2,3) (2,)
```

What happens here is that Numpy, again, adds a dimension to `b`, making it of shape (1, 2). The sizes of the last dimensions of `b` and `c` (2 and 3, respectively) are then compared and found to differ. Since none of these dimensions is of size 1 (therefore, unstretchable) Numpy gives up and produces an error.

The solution to multiplying `c` and `b` above is to specifically tell Numpy that it must add that extra dimension as the second dimension of `b`. This is done by using `None` to index that second dimension. The shape of `b` then becomes (2, 1), which is compatible for broadcasting with `c`:

```
>>> c = np.array([
    [0, 1, 2],
    [3, 4, 5],
    ])
>>> b = np.array([10, 100])
>>> c * b[:, None]
array([[ 0, 10, 20],
       [300, 400, 500]])
```

A good visual description of these rules, together with some advanced broadcasting applications can be found in this [tutorial of Numpy broadcasting rules](#).

## 6.5 dtype

A commonly used term in working with numpy is `dtype` - short for data type. This is typically `int` or `float`, followed by some number, e.g. `int8`. This means the value is integer with a size of 8 bits. As an example, let's discuss the properties of an `int8`.

Each bit is either 0 or 1. With 8 of them, we have  $2^8 = 256$  possible values. Since we also have to count zero itself, the largest possible value is 255. The data type we have now described is called `uint8`, where the `u` stands for unsigned: only positive values are allowed. If we want to allow negative numbers we use `int8`. The range then shifts to -128 to +127.

The same holds for bigger numbers. An `int64` for example is a 64 bit unsigned integer with a range of -9223372036854775808 to 9223372036854775807. It is also the standard type on a 64 bits machine. You might think bigger is better. You'd be wrong. If you know the elements of your array are never going to be bigger than 100, why waste all the memory space? You might be better off setting your array to `uint8` to conserve memory. In general however, the default setting is fine. Only when you run into memory related problems should you remember this comment.

What happens when you set numbers bigger than the maximum value of your dtype?

```
>>> import numpy as np
>>> a = np.array([200], dtype='uint8')
>>> a + a
array([144], dtype=uint8)
```

That doesn't seem right, does it? If you add two `uint8`, the result of  $200 + 200$  cannot be 400, because that doesn't fit in a `uint8`. In standard Python, Python does a lot of magic in the background to make sure the result is the 400 you would expect. But numpy doesn't, and will return 144. Why 144 is left as an exercise. To fix this, you should make sure that your numbers were not stored as `uint8`, but as something larger; `uint16` for example. That way the resulting 400 will fit.

```
>>> import numpy as np
>>> a = np.array([200], dtype='uint16')
>>> a + a
array([400], dtype=uint16)
```

By now you must be thinking: so bigger is better after all! Just use the biggest possible int all the time, and you'll be fine! Apart from the fact that there is no biggest int, there is a bigger problem. If you work with images, each pixel from that image is stored as an RGB tuple: the intensity in red, green and blue. Each of these is a `uint8` value for most standard formats such as .jpg and .png. For example, (0, 0, 0) will be black, and (255, 0, 0) is red. This means

that when you load an image from your hard drive this dtype is selected for you, and if you are not aware of this, what will happen when you add an image to itself? (In other words, place two copies on top of each other) You might expect that everything will become more dense. Instead, you'll get noise because of the effect we just talked about.

## 6.6 Changing dtype

To change the dtype of an existing array, you can use the `astype` method:

```
>>> import numpy as np
>>> a = np.array([200], dtype='uint8')
>>> a.astype('uint64')
```

## 6.7 Advanced Usage

Numpy has vast capabilities. It has way too many options to discuss here. More information can be found in

1. the [Quickstart Numpy Tutorial](#);
2. the [Numpy indexing documentation](#) (for advanced slicing and indexing);
3. and the [Numpy broadcasting rules](#) (for what happens when performing operations between arrays of different shapes and sizes).

## 6.8 Exercises

1. Make an array with `dtype = uint8` and elements of your choosing. Keep adding to it until (one of) the items go over 255. What happens? Hint: make an array, and just add a constant to it. The constant will be added to all the items of the array element-wise.
2. Use a mask to multiply all values below 100 in the following list by 2:

```
>>> a = np.array([230, 10, 284, 39, 76])
```

Repeat this until all values are above 100. (Not manually, but by looping)

Then, select all values between  $150 < a < 200$ .

## 7.1 About files

While a program is running, its data is stored in *random access memory* (RAM). RAM is fast and inexpensive, but it is also **volatile**, which means that when the program ends, or the computer shuts down, data in RAM disappears. To make data available the next time the computer is turned on and the program is started, it has to be written to a **non-volatile** storage medium, such a hard drive, usb drive, or CD-RW.

Data on non-volatile storage media is stored in named locations called **files**. By reading and writing files, programs can save information between program runs.

Working with files is a lot like working with a notebook. To use a notebook, it has to be opened. When done, it has to be closed. While the notebook is open, it can either be read from or written to. In either case, the notebook holder knows where they are. They can read the whole notebook in its natural order or they can skip around.

All of this applies to files as well. To open a file, we specify its name and indicate whether we want to read or write.

## 7.2 Writing our first file

Let's begin with a simple program that writes three lines of text into a file:

```
1 with open("test.txt", "w") as myfile:
2     myfile.write("My first file written from Python\n")
3     myfile.write("-----\n")
4     myfile.write("Hello, world!\n")
```

Opening a file creates what we call a file **handle**. In this example, the variable `myfile` refers to the new handle object. Our program calls methods on the handle, and this makes changes to the actual file which is usually located on our disk.

On line 1, the `open` function takes two arguments. The first is the name of the file, and the second is the **mode**. Mode `"w"` means that we are opening the file for writing.

With mode "w", if there is no file named `test.txt` on the disk, it will be created. If there already is one, it will be replaced by the file we are writing.

To put data in the file we invoke the `write` method on the handle, shown in lines 2, 3 and 4 above. In bigger programs, lines 2–4 will usually be replaced by a loop that writes many more lines into the file.

The file is closed after line 4, at the end of the `with` block. A `with` block make sure that the file get close even if an error occurs (power outages excluded).

---

### A handle is somewhat like a TV remote control

We're all familiar with a remote control for a TV. We perform operations on the remote control — switch channels, change the volume, etc. But the real action happens on the TV. So, by simple analogy, we'd call the remote control our *handle* to the underlying TV.

Sometimes we want to emphasize the difference — the file handle is not the same as the file, and the remote control is not the same as the TV. But at other times we prefer to treat them as a single mental chunk, or abstraction, and we'll just say "close the file", or "flip the TV channel".

---

## 7.3 Reading a file line-at-a-time

Now that the file exists on our disk, we can open it, this time for reading, and read all the lines in the file, one at a time. This time, the mode argument is "r" for reading:

```
1 with open("test.txt", "r") as my_new_handle:
2     for the_line in my_new_handle:
3         # Do something with the line we just read.
4         # Here we just print it.
5         print(the_line, end="")
```

This is a handy pattern for our toolbox. In bigger programs, we'd squeeze more extensive logic into the body of the loop at line 5 — for example, if each line of the file contained the name and email address of one of our friends, perhaps we'd split the line into some pieces and call a function to send the friend a party invitation.

On line 5 we suppress the newline character that `print` usually appends to our strings with `end=""`. Why? This is because the string already has its own newline: the `for` statement in line 2 reads everything up to *and including* the newline character.

If we try to open a file that doesn't exist, we get an error:

```
>>> mynewhandle = open("wharrah.txt", "r")
FileNotFoundError: [Errno 2] No such file or directory: "wharrah.txt"
```

## 7.4 Turning a file into a list of lines

It is often useful to fetch data from a disk file and turn it into a list of lines. Suppose we have a file containing our friends and their email addresses, one per line in the file. But we'd like the lines sorted into alphabetical order. A good plan is to read everything into a list of lines, then sort the list, and then write the sorted list back to another file:

```
1 with open("friends.txt", "r") as input_file:
2     all_lines = input_file.readlines()
3
```

(continues on next page)

(continued from previous page)

```
4 all_lines.sort()
5
6 with open("sortedfriends.txt", "w") as output_file:
7     for line in all_lines:
8         outut_file.write(line)
```

The `readlines` method in line 2 reads all the lines and returns a list of the strings.

We could have used the template from the previous section to read each line one-at-a-time, and to build up the list ourselves, but it is a lot easier to use the method that the Python implementors gave us!

## 7.5 Reading the whole file at once

Another way of working with text files is to read the complete contents of the file into a string, and then to use our string-processing skills to work with the contents.

We'd normally use this method of processing files if we were not interested in the line structure of the file. For example, we've seen the `split` method on strings which can break a string into words. So here is how we might count the number of words in a file:

```
1 with open("somefile.txt") as f:
2     content = f.read()
3     words = content.split()
4     print("There are {0} words in the file.".format(len(words)))
```

Notice here that we left out the `"r"` mode in line 1. By default, if we don't supply the mode, Python opens the file for reading.

---

### Your file paths may need to be explicitly named.

In the above example, we're assuming that the file `somefile.txt` is in the same directory as your Python source code. If this is not the case, you may need to provide a full or a relative path to the file. On Windows, a full path could look like `"C:\\temp\\somefile.txt"`, while on a Unix system the full path could be `"/home/jimmy/somefile.txt"`.

We'll return to this later in this chapter.

---

## 7.6 An example

Many useful line-processing programs will read a text file line-at-a-time and do some minor processing as they write the lines to an output file. They might number the lines in the output file, or insert extra blank lines after every 60 lines to make it convenient for printing on sheets of paper, or extract some specific columns only from each line in the source file, or only print lines that contain a specific substring. We call this kind of program a **filter**.

Here is a filter that copies one file to another, omitting any lines that begin with `#`:

```
1 def filter(oldfile, newfile):
2     with open(oldfile, "r") as infile, open(newfile, "w") as outfile:
3         for line in infile:
4             # Put any processing logic here
5             if not line.startswith('#'):
6                 outfile.write(line)
```



On line 2, we open two files: the file to read, and the file to write. From line 3, we read the input file line by line. We write the line in the output file only if the condition on line 5 is true.

## 7.7 Directories

Files on non-volatile storage media are organized by a set of rules known as a **file system**. File systems are made up of files and **directories**, which are containers for both files and other directories.

When we create a new file by opening it and writing, the new file goes in the current directory (wherever we were when we ran the program). Similarly, when we open a file for reading, Python looks for it in the current directory.

If we want to open a file somewhere else, we have to specify the **path** to the file, which is the name of the directory (or folder) where the file is located:

```
>>> wordsfile = open("/usr/share/dict/words", "r")
>>> wordlist = wordsfile.readlines()
>>> print(wordlist[:6])
['\n', 'A\n', 'A's\n', 'AOL\n', 'AOL's\n', 'Aachen\n']
```

This (Unix) example opens a file named `words` that resides in a directory named `dict`, which resides in `share`, which resides in `usr`, which resides in the top-level directory of the system, called `/`. It then reads in each line into a list using `readlines`, and prints out the first 5 elements from that list.

A Windows path might be `"c:/temp/words.txt"` or `"c:\\temp\\words.txt"`. Because backslashes are used to escape things like newlines and tabs, we need to write two backslashes in a literal string to get one! So the length of these two strings is the same!

We cannot use `/` or `\` as part of a filename; they are reserved as a **delimiter** between directory and filenames.

The file `/usr/share/dict/words` should exist on Unix-based systems, and contains a list of words in alphabetical order.

## 7.8 What about fetching something from the web?

The Python libraries are pretty messy in places. But here is a very simple example that copies the contents at some web URL to a local file.

```
1 import urllib.request
2
3 url = "http://xml.resource.org/public/rfc/txt/rfc793.txt"
4 destination_filename = "rfc793.txt"
5
6 urllib.request.urlretrieve(url, destination_filename)
```

The `urlretrieve` function — just one call — could be used to download any kind of content from the Internet.

**We'll need to get a few things right before this works:**

- The resource we're trying to fetch must exist! Check this using a browser.
- We'll need permission to write to the destination filename, and the file will be created in the "current directory" - i.e. the same folder that the Python program is saved in.
- If we are behind a proxy server that requires authentication, (as some students are), this may require some more special handling to work around our proxy. Use a local resource for the purpose of this demonstration!

Here is a slightly different example using the *requests* module. This module is not part of the standard library distributed with python, however it is easier to use and significantly more potent than the *urllib* module distributed with python. Read *requests* documentation on <http://docs.python-requests.org> to learn how to install and use the module. Here, rather than save the web resource to our local disk, we read it directly into a string, and we print that string:

```
1 import requests
2
3 url = "http://xml.resource.org/public/rfc/txt/rfc793.txt"
4 response = requests.get(url)
5 print(response.text)
```

Opening the remote URL returns the response from the server. That response contains several types of information, and the *requests* module allows us to access them in various ways. On line 5, we get the downloaded document as a single string. We could also read it line by line as follows:

```
1 import requests
2
3 url = "http://xml.resource.org/public/rfc/txt/rfc793.txt"
4 response = requests.get(url)
5 for line in response:
6     print(line)
```

## 7.9 Glossary

**delimiter** A sequence of one or more characters used to specify the boundary between separate parts of text.

**directory** A named collection of files, also called a folder. Directories can contain files and other directories, which are referred to as *subdirectories* of the directory that contains them.

**file** A named entity, usually stored on a hard drive, floppy disk, or CD-ROM, that contains a stream of characters.

**file system** A method for naming, accessing, and organizing files and the data they contain.

**handle** An object in our program that is connected to an underlying resource (e.g. a file). The file handle lets our program manipulate/read/write/close the actual file that is on our disk.

**mode** A distinct method of operation within a computer program. Files in Python can be opened in one of four modes: read ("r"), write ("w"), append ("a"), and read and write ("+").

**non-volatile memory** Memory that can maintain its state without power. Hard drives, flash drives, and rewritable compact disks (CD-RW) are each examples of non-volatile memory.

**path** A sequence of directory names that specifies the exact location of a file.

**text file** A file that contains printable characters organized into lines separated by newline characters.

**volatile memory** Memory which requires an electrical current to maintain state. The *main memory* or RAM of a computer is volatile. Information stored in RAM is lost when the computer is turned off.

## 7.10 Exercises

1. Write a program that reads a file and writes out a new file with the lines in reversed order (i.e. the first line in the old file becomes the last one in the new file.)
2. Write a program that reads a file and prints only those lines that contain the substring *snake*.

3. Write a program that reads a text file and produces an output file which is a copy of the file, except the first five columns of each line contain a four digit line number, followed by a space. Start numbering the first line in the output file at 1. Ensure that every line number is formatted to the same width in the output file. Use one of your Python programs as test data for this exercise: your output should be a printed and numbered listing of the Python program.
4. Write a program that undoes the numbering of the previous exercise: it should read a file with numbered lines and produce another file without line numbers.
5. Write a program that takes the dictionary used above, and returns some of the words using 1337sp34k

A **module** is a file containing Python definitions and statements intended for use in other Python programs. There are many Python modules that come with Python as part of the **standard library**. We have seen at least two of these already, the `turtle` module and the `string` module.

We have also shown you how to access help. The help system contains a listing of all the standard modules that are available with Python. Play with help!

## 8.1 Random numbers

We often want to use random numbers in programs, here are a few typical uses:

- To play a game of chance where the computer needs to throw some dice, pick a number, or flip a coin,
- To shuffle a deck of playing cards randomly,
- To allow/make an enemy spaceship appear at a random location and start shooting at the player,
- To simulate possible rainfall when we make a computerized model for estimating the environmental impact of building a dam,
- For encrypting banking sessions on the Internet.

Python provides a module `random` that helps with tasks like this. You can look it up using help, but here are the key things we'll do with it:

```
1 import random
2
3 # Create a black box object that generates random numbers
4 rng = random.Random()
5
6 dice_throw = rng.randrange(1,7)    # Return an int, one of 1,2,3,4,5,6
7 delay_in_seconds = rng.random() * 5.0
```

The `randrange` method call generates an integer between its lower and upper argument, using the same semantics as `range` — so the lower bound is included, but the upper bound is excluded. All the values have an equal probability

of occurring (i.e. the results are *uniformly* distributed). Like `range`, `randrange` can also take an optional step argument. So let's assume we needed a random odd number less than 100, we could say:

```
1 random_odd = rng.randrange(1, 100, 2)
```

Other methods can also generate other distributions e.g. a bell-shaped, or “normal” distribution might be more appropriate for estimating seasonal rainfall, or the concentration of a compound in the body after taking a dose of medicine.

The `random` method returns a floating point number in the interval `[0.0, 1.0)` — the square bracket means “closed interval on the left” and the round parenthesis means “open interval on the right”. In other words, 0.0 is possible, but all returned numbers will be strictly less than 1.0. It is usual to *scale* the results after calling this method, to get them into an interval suitable for your application. In the case shown here, we've converted the result of the method call to a number in the interval `[0.0, 5.0)`. Once more, these are uniformly distributed numbers — numbers close to 0 are just as likely to occur as numbers close to 0.5, or numbers close to 1.0.

This example shows how to shuffle a list. (`shuffle` cannot work directly with a lazy promise, so notice that we had to convert the range object using the `list` type converter first.)

```
1 cards = list(range(52)) # Generate ints [0 .. 51]
2                               # representing a pack of cards.
3 rng.shuffle(cards)      # Shuffle the pack
```

### 8.1.1 Repeatability and Testing

Random number generators are based on a **deterministic** algorithm — repeatable and predictable. So they're called **pseudo-random** generators — they are not genuinely random. They start with a *seed* value. Each time you ask for another random number, you'll get one based on the current seed attribute, and the state of the seed (which is one of the attributes of the generator) will be updated.

For debugging and for writing unit tests, it is convenient to have repeatability — programs that do the same thing every time they are run. We can arrange this by forcing the random number generator to be initialized with a known seed every time. (Often this is only wanted during testing — playing a game of cards where the shuffled deck was always in the same order as last time you played would get boring very rapidly!)

```
1 drng = random.Random(123) # Create generator with known starting state
```

This alternative way of creating a random number generator gives an explicit seed value to the object. Without this argument, the system probably uses something based on the time. So grabbing some random numbers from `drng` today will give you precisely the same random sequence as it will tomorrow!

### 8.1.2 Picking balls from bags, throwing dice, shuffling a pack of cards

Here is an example to generate a list containing  $n$  random ints between a lower and an upper bound:

```
1 import random
2
3 def make_random_ints(num, lower_bound, upper_bound):
4     """
5     Generate a list containing num random ints between lower_bound
6     and upper_bound. upper_bound is an open bound.
7     """
8     rng = random.Random() # Create a random number generator
9     result = []
10    for i in range(num):
```

(continues on next page)

(continued from previous page)

```
11     result.append(rng.randrange(lower_bound, upper_bound))
12     return result
```

```
>>> make_random_ints(5, 1, 13)  # Pick 5 random month numbers
[8, 1, 8, 5, 6]
```

Notice that we got a duplicate in the result. Often this is wanted, e.g. if we throw a die five times, we would expect some duplicates.

But what if you don't want duplicates? If you wanted 5 distinct months, then this algorithm is wrong. In this case a good algorithm is to generate the list of possibilities, shuffle it, and slice off the number of elements you want:

```
1  xs = list(range(1,13))  # Make list 1..12 (there are no duplicates)
2  rng = random.Random()   # Make a random number generator
3  rng.shuffle(xs)         # Shuffle the list
4  result = xs[:5]         # Take the first five elements
```

In statistics courses, the first case — allowing duplicates — is usually described as pulling balls out of a bag *with replacement* — you put the drawn ball back in each time, so it can occur again. The latter case, with no duplicates, is usually described as pulling balls out of the bag *without replacement*. Once the ball is drawn, it doesn't go back to be drawn again. TV lotto games work like this.

The second “shuffle and slice” algorithm would not be so great if you only wanted a few elements, but from a very large domain. Suppose I wanted five numbers between one and ten million, without duplicates. Generating a list of ten million items, shuffling it, and then slicing off the first five would be a performance disaster! So let us have another try:

```
1  import random
2
3  def make_random_ints_no_dups(num, lower_bound, upper_bound):
4      """
5      Generate a list containing num random ints between
6      lower_bound and upper_bound. upper_bound is an open bound.
7      The result list cannot contain duplicates.
8      """
9      result = []
10     rng = random.Random()
11     for i in range(num):
12         while True:
13             candidate = rng.randrange(lower_bound, upper_bound)
14             if candidate not in result:
15                 break
16             result.append(candidate)
17     return result
18
19 xs = make_random_ints_no_dups(5, 1, 10000000)
20 print(xs)
```

This agreeably produces 5 random numbers, without duplicates:

```
[3344629, 1735163, 9433892, 1081511, 4923270]
```

Even this function has its pitfalls. Can you spot what is going to happen in this case?

```
1  xs = make_random_ints_no_dups(10, 1, 6)
```

## 8.2 The `time` module

As we start to work with more sophisticated algorithms and bigger programs, a natural concern is “*is our code efficient?*” One way to experiment is to time how long various operations take. The `time` module has a function called `clock` that is recommended for this purpose. Whenever `clock` is called, it returns a floating point number representing how many seconds have elapsed since your program started running.

The way to use it is to call `clock` and assign the result to a variable, say `t0`, just before you start executing the code you want to measure. Then after execution, call `clock` again, (this time we’ll save the result in variable `t1`). The difference `t1-t0` is the time elapsed, and is a measure of how fast your program is running.

Let’s try a small example. Python has a built-in `sum` function that can sum the elements in a list. We can also write our own. How do we think they would compare for speed? We’ll try to do the summation of a list `[0, 1, 2 ...]` in both cases, and compare the results:

```
1 import time
2
3 def do_my_sum(xs):
4     sum = 0
5     for v in xs:
6         sum += v
7     return sum
8
9 sz = 10000000          # Lets have 10 million elements in the list
10 testdata = range(sz)
11
12 t0 = time.clock()
13 my_result = do_my_sum(testdata)
14 t1 = time.clock()
15 print("my_result      = {0} (time taken = {1:.4f} seconds)"
16       .format(my_result, t1-t0))
17
18 t2 = time.clock()
19 their_result = sum(testdata)
20 t3 = time.clock()
21 print("their_result = {0} (time taken = {1:.4f} seconds)"
22       .format(their_result, t3-t2))
```

On a reasonably modest laptop, we get these results:

```
my_sum      = 49999995000000 (time taken = 1.5567 seconds)
their_sum   = 49999995000000 (time taken = 0.9897 seconds)
```

So our function runs about 57% slower than the built-in one. Generating and summing up ten million elements in under a second is not too shabby!

## 8.3 The `math` module

The `math` module contains the kinds of mathematical functions you’d typically find on your calculator (`sin`, `cos`, `sqrt`, `asin`, `log`, `log10`) and some mathematical constants like `pi` and `e`:

```
>>> import math
>>> math.pi          # Constant pi
3.141592653589793
>>> math.e           # Constant natural log base
```

(continues on next page)

(continued from previous page)

```
2.718281828459045
>>> math.sqrt(2.0)           # Square root function
1.4142135623730951
>>> math.radians(90)         # Convert 90 degrees to radians
1.5707963267948966
>>> math.sin(math.radians(90)) # Find sin of 90 degrees
1.0
>>> math.asin(1.0) * 2       # Double the arcsin of 1.0 to get pi
3.141592653589793
```

Like almost all other programming languages, angles are expressed in *radians* rather than degrees. There are two functions `radians` and `degrees` to convert between these two popular ways of measuring angles.

Notice another difference between this module and our use of `random` and `turtle`: in `random` and `turtle` we create objects and we call methods on the object. This is because objects have *state* — a turtle has a color, a position, a heading, etc., and every random number generator has a seed value that determines its next result.

Mathematical functions are “pure” and don’t have any state — calculating the square root of 2.0 doesn’t depend on any kind of state or history about what happened in the past. So the functions are not methods of an object — they are simply functions that are grouped together in a module called `math`.

## 8.4 Creating your own modules

All we need to do to create our own modules is to save our script as a file with a `.py` extension. Suppose, for example, this script is saved as a file named `seqtools.py`:

```
1 def remove_at(pos, seq):
2     return seq[:pos] + seq[pos+1:]
```

We can now use our module, both in scripts we write, or in the interactive Python interpreter. To do so, we must first `import` the module.

```
>>> import seqtools
>>> s = "A string!"
>>> seqtools.remove_at(4, s)
'A sting!'
```

We do not include the `.py` file extension when importing. Python expects the file names of Python modules to end in `.py`, so the file extension is not included in the **import statement**.

The use of modules makes it possible to break up very large programs into manageable sized parts, and to keep related parts together.

## 8.5 Namespaces

A **namespace** is a collection of identifiers that belong to a module, or to a function, (and as we will see soon, in classes too). Generally, we like a namespace to hold “related” things, e.g. all the math functions, or all the typical things we’d do with random numbers.

Each module has its own namespace, so we can use the same identifier name in multiple modules without causing an identification problem.



```
1 # module1.py
2
3 question = "What is the meaning of Life, the Universe, and Everything?"
4 answer = 42
```

```
1 # module2.py
2
3 question = "What is your quest?"
4 answer = "To seek the holy grail."
```

We can now import both modules and access `question` and `answer` in each:

```
1 import module1
2 import module2
3
4 print(module1.question)
5 print(module2.question)
6 print(module1.answer)
7 print(module2.answer)
```

will output the following:

```
What is the meaning of Life, the Universe, and Everything?
What is your quest?
42
To seek the holy grail.
```

Functions also have their own namespaces:

```
1 def f():
2     n = 7
3     print("printing n inside of f:", n)
4
5 def g():
6     n = 42
7     print("printing n inside of g:", n)
8
9 n = 11
10 print("printing n before calling f:", n)
11 f()
12 print("printing n after calling f:", n)
13 g()
14 print("printing n after calling g:", n)
```

Running this program produces the following output:

```
printing n before calling f: 11
printing n inside of f: 7
printing n after calling f: 11
printing n inside of g: 42
printing n after calling g: 11
```

The three `n`'s here do not collide since they are each in a different namespace — they are three names for three different variables, just like there might be three different instances of people, all called “Bruce”.

Namespaces permit several programmers to work on the same project without having naming collisions.

---

### How are namespaces, files and modules related?

Python has a convenient and simplifying one-to-one mapping, one module per file, giving rise to one namespace. Also, Python takes the module name from the file name, and this becomes the name of the namespace. `math.py` is a filename, the module is called `math`, and its namespace is `math`. So in Python the concepts are more or less interchangeable.

But you will encounter other languages (e.g. C#), that allow one module to span multiple files, or one file to have multiple namespaces, or many files to all share the same namespace. So the name of the file doesn't need to be the same as the namespace.

So a good idea is to try to keep the concepts distinct in your mind.

Files and directories organize *where* things are stored in our computer. On the other hand, namespaces and modules are a programming concept: they help us organize how we want to group related functions and attributes. They are not about “where” to store things, and should not have to coincide with the file and directory structures.

So in Python, if you rename the file `math.py`, its module name also changes, your `import` statements would need to change, and your code that refers to functions or attributes inside that namespace would also need to change.

In other languages this is not necessarily the case. So don't blur the concepts, just because Python blurs them!

---

## 8.6 Scope and lookup rules

The **scope** of an identifier is the region of program code in which the identifier can be accessed, or used.

There are three important scopes in Python:

- **Local scope** refers to identifiers declared within a function. These identifiers are kept in the namespace that belongs to the function, and each function has its own namespace.
- **Global scope** refers to all the identifiers declared within the current module, or file.
- **Built-in scope** refers to all the identifiers built into Python — those like `range` and `min` that can be used without having to import anything, and are (almost) always available.

Python can help you by telling you what is in which scope. Use the functions `locals`, `globals`, and `dir` to see for yourself!

Python (like most other computer languages) uses precedence rules: the same name could occur in more than one of these scopes, but the innermost, or local scope, will always take precedence over the global scope, and the global scope always gets used in preference to the built-in scope. Let's start with a simple example:

```
1 def range(n):
2     return 123*n
3
4 print(range(10))
```

What gets printed? We've defined our own function called `range`, so there is now a potential ambiguity. When we use `range`, do we mean our own one, or the built-in one? Using the scope lookup rules determines this: our own `range` function, not the built-in one, is called, because our function `range` is in the global namespace, which takes precedence over the built-in names.

So although names like `range` and `min` are built-in, they can be “hidden” from your use if you choose to define your own variables or functions that reuse those names. (It is a confusing practice to redefine built-in names — so to be a good programmer you need to understand the scope rules and understand that you can do nasty things that will cause confusion, and then you avoid doing them!)

Now, a slightly more complex example:

```
1  n = 10
2  m = 3
3  def f(n):
4      m = 7
5      return 2*n+m
6
7  print(f(5), n, m)
```

This prints `17 10 3`. The reason is that the two variables `m` and `n` in lines 1 and 2 are outside the function in the global namespace. Inside the function, new variables called `n` and `m` are created *just for the duration of the execution of f*. These are created in the local namespace of function `f`. Within the body of `f`, the scope lookup rules determine that we use the local variables `m` and `n`. By contrast, after we’ve returned from `f`, the `n` and `m` arguments to the `print` function refer to the original variables on lines 1 and 2, and these have not been changed in any way by executing function `f`.

Notice too that the `def` puts name `f` into the global namespace here. So it can be called on line 7.

What is the scope of the variable `n` on line 1? Its scope — the region in which it is visible — is lines 1, 2, 6, 7. It is hidden from view in lines 3, 4, 5 because of the local variable `n`.

## 8.7 Attributes and the dot operator

Variables defined inside a module are called **attributes** of the module. We’ve seen that objects have attributes too: for example, most objects have a `__doc__` attribute, some functions have a `__annotations__` attribute. Attributes are accessed using the **dot operator** (`.`). The `question` attribute of `module1` and `module2` is accessed using `module1.question` and `module2.question`.

Modules contain functions as well as attributes, and the dot operator is used to access them in the same way. `seqtools.remove_at` refers to the `remove_at` function in the `seqtools` module.

When we use a dotted name, we often refer to it as a **fully qualified name**, because we’re saying exactly which `question` attribute we mean.

## 8.8 Three `import` statement variants

Here are three different ways to import names into the current namespace, and to use them:

```
1  import math
2  x = math.sqrt(10)
```

Here just the single identifier `math` is added to the current namespace. If you want to access one of the functions in the module, you need to use the dot notation to get to it.

Here is a different arrangement:

```
1  from math import cos, sin, sqrt
2  x = sqrt(10)
```

The names are added directly to the current namespace, and can be used without qualification. The name `math` is not itself imported, so trying to use the qualified form `math.sqrt` would give an error.

Then we have a convenient shorthand:

```
1  from math import *      # Import all the identifiers from math,
2                          #   adding them to the current namespace.
3  x = sqrt(10)           # Use them without qualification.
```

Of these three, the first method is generally preferred, even though it means a little more typing each time. Although, we can make things shorter by importing a module under a different name:

```
1  >>> import math as m
2  >>> m.pi
3  3.141592653589793
```

But hey, with nice editors that do auto-completion, and fast fingers, that's a small price!

Finally, observe this case:

```
1  def area(radius):
2      import math
3      return math.pi * radius * radius
4
5  x = math.sqrt(10)      # This gives an error
```

Here we imported `math`, but we imported it into the local namespace of `area`. So the name is usable within the function body, but not in the enclosing script, because it is not in the global namespace.

## 8.9 Glossary

**attribute** A variable defined inside a module (or class or instance – as we will see later). Module attributes are accessed by using the **dot operator** (`.`).

**dot operator** The dot operator (`.`) permits access to attributes and functions of a module (or attributes and methods of a class or instance – as we have seen elsewhere).

**fully qualified name** A name that is prefixed by some namespace identifier and the dot operator, or by an instance object, e.g. `math.sqrt` or `tess.forward(10)`.

**import statement** A statement which makes the objects contained in a module available for use within another module. There are two forms for the import statement. Using hypothetical modules named `mymod1` and `mymod2` each containing functions `f1` and `f2`, and variables `v1` and `v2`, examples of these two forms include:

```
1  import mymod1
2  from mymod2 import f1, f2, v1, v2
```

The second form brings the imported objects into the namespace of the importing module, while the first form preserves a separate namespace for the imported module, requiring `mymod1.v1` to access the `v1` variable from that module.

**method** Function-like attribute of an object. Methods are *invoked* (called) on an object using the dot operator. For example:

```
>>> s = "this is a string."
>>> s.upper()
```

(continues on next page)

(continued from previous page)

```
'THIS IS A STRING.'  
>>>
```

We say that the method, `upper` is invoked on the string, `s`. `s` is implicitly the first argument to `upper`.

**module** A file containing Python definitions and statements intended for use in other Python programs. The contents of a module are made available to the other program by using the `import` statement.

**namespace** A syntactic container providing a context for names so that the same name can reside in different namespaces without ambiguity. In Python, modules, classes, functions and methods all form namespaces.

**naming collision** A situation in which two or more names in a given namespace cannot be unambiguously resolved. Using

```
1 import string
```

instead of

```
1 from string import *
```

prevents naming collisions.

**standard library** A library is a collection of software used as tools in the development of other software. The standard library of a programming language is the set of such tools that are distributed with the core programming language. Python comes with an extensive standard library.

## 8.10 Exercises

1. Open help for the `calendar` module.

- (a) Try the following:

```
1 import calendar  
2 cal = calendar.TextCalendar()           # Create an instance  
3 cal.pryear(2012)                        # What happens here?
```

- (b) Observe that the week starts on Monday. An adventurous CompSci student believes that it is better mental chunking to have his week start on Thursday, because then there are only two working days to the weekend, and every week has a break in the middle. Read the documentation for `TextCalendar`, and see how you can help him print a calendar that suits his needs.
- (c) Find a function to print just the month in which your birthday occurs this year.
- (d) Try this:

```
1 d = calendar.LocaleTextCalendar(6, "SPANISH")  
2 d.pryear(2012)
```

Try a few other languages, including one that doesn't work, and see what happens.

- (e) Experiment with `calendar.isleap`. What does it expect as an argument? What does it return as a result? What kind of a function is this?

Make detailed notes about what you learned from these exercises.

2. Open help for the `math` module.

- (a) How many functions are in the `math` module?

- (b) What does `math.ceil` do? What about `math.floor`? (*hint: both floor and ceil expect floating point arguments.*)
- (c) Describe how we have been computing the same value as `math.sqrt` without using the `math` module.
- (d) What are the two data constants in the `math` module?

Record detailed notes of your investigation in this exercise.

3. Investigate the `copy` module. What does `deepcopy` do? In which exercises from last chapter would `deepcopy` have come in handy?
4. Create a module named `mymodule1.py`. Add attributes `myage` set to your current age, and `year` set to the current year. Create another module named `mymodule2.py`. Add attributes `myage` set to 0, and `year` set to the year you were born. Now create a file named `namespace_test.py`. Import both of the modules above and write the following statement:

```
1 print( (mymodule2.myage - mymodule1.myage) ==
2        (mymodule2.year - mymodule1.year) )
```

When you will run `namespace_test.py` you will see either `True` or `False` as output depending on whether or not you've already had your birthday this year.

What this example illustrates is that out different modules can both have attributes named `myage` and `year`. Because they're in different namespaces, they don't clash with one another. When we write `namespace_test.py`, we fully qualify exactly which variable `year` or `myage` we are referring to.

5. Add the following statement to `mymodule1.py`, `mymodule2.py`, and `namespace_test.py` from the previous exercise:

```
1 print("My name is", __name__)
```

Run `namespace_test.py`. What happens? Why? Now add the following to the bottom of `mymodule1.py`:

```
1 if __name__ == "__main__":
2     print("This won't run if I'm imported.")
```

Run `mymodule1.py` and `namespace_test.py` again. In which case do you see the new print statement?

6. In a Python shell / interactive interpreter, try the following:

```
>>> import this
```

What does Tim Peters have to say about namespaces?



## CHAPTER 9

---

### More datatypes

---

You have already encountered the most important datatypes Python has to offer: booleans, ints, floats, strings, tuples, lists and dictionaries. However, there is more to it than hinted at previously. In this section, we will focus mainly on tuples and lists, and introduce sets and frozensets.

### 9.1 Mutable versus immutable and aliasing

Some datatypes in Python are **mutable**. This means their contents can be changed after they have been created. Lists and dictionaries are good examples of mutable datatypes.

```
>>> my_list = [2, 4, 5, 3, 6, 1]
>>> my_list[0] = 9
>>> my_list
[9, 4, 5, 3, 6, 1]
```

Tuples and strings are examples of immutable datatypes, their contents can not be changed after they have been created:

```
>>> my_tuple = (2, 5, 3, 1)
>>> my_tuple[0] = 9
Traceback (most recent call last):
  File "<interactive input>", line 2, in <module>
TypeError: 'tuple' object does not support item assignment
>>>
```

Mutability is usually useful, but it may lead to something called aliasing. In this case, two variables refer to the same object and mutating one will also change the other:

```
>>> list_one = [1, 2, 3, 4, 6]
>>> list_two = list_one
>>> list_two[-1] = 5
>>> list_one
[1, 2, 3, 4, 5]
```



This happens, because both `list_one` and `list_two` refer to the same memory address containing the actual list. You can check this using the built-in function `id`:

```
>>> list_one = [1, 2, 3, 4, 6]
>>> list_two = list_one
>>> id(list_one) == id(list_two)
True
```

You can escape this problem by making a copy of the list:

```
>>> list_one = [1, 2, 3, 4, 6]
>>> list_two = list_one[:]
>>> id(list_one) == id(list_two)
False
>>> list_two[-1] = 5
>>> list_two
[1, 2, 3, 4, 5]
>>> list_one
[1, 2, 3, 4, 6]
```

However, this will not work for nested lists because of the same reason. The module `copy` provides functions to solve this.

## 9.2 Sets and frozensets

Given that tuples and lists are ordered, and dictionaries are unordered, we can construct the following table.

	Ordered	Unordered
Mutable	list	dict
Immutable	tuple	

This reveals an empty spot: we don't know any immutable, unordered datatypes yet. Additionally, you can argue that a dictionary doesn't belong in this table, since it is a *mapping* type whilst lists and tuples are not: a dictionary maps keys to values. This is where sets and frozensets come in. A **set** is an unordered, mutable datatype; and a **frozenset** is an unordered, immutable datatype.

	Ordered	Unordered
Mutable	list	set
Immutable	tuple	frozenset

Since sets and frozensets are unordered, they share some properties with dictionaries: for example, it's elements are unique. Creating a set, and adding elements to it is simple.

```
>>> my_set = set([1, 4, 2, 3, 4])
>>> my_set
{1, 2, 3, 4}
>>> my_set.add(13)
>>> my_set
{1, 2, 3, 4, 13}
```

Sets may seem sorted in the example above, but this is completely coincidental. Sets also support common operations such as membership testing (`3 in my_set`); and iteration (`for x in my_set:`). Additionally, you can add and subtract sets from each other:

```
1 set1 = set([1, 2, 3])
2 set2 = set([4, 5, 6])
3 print(set1 | set2)  # {1, 2, 3, 4, 5, 6}
4 print(set1 & set2)  # set()
5 set2 = set([2, 3, 4, 5])
6 print(set1 & set2)  # {2, 3}
7 print(set1 - set2)  # {1}
```

Frozensets are mostly the same as set, other than that they can not be modified; i.e. you can't add or remove items. See also the documentation [online](#).

More exotic data types - such as queues, stacks and ordered dictionaries - are provided in Python's `collections` module. You can find the documentation [here](#).



**Recursion** means “defining something in terms of itself” usually at some smaller scale, perhaps multiple times, to achieve your objective. For example, we might say “A human being is someone whose mother is a human being”, or “a directory is a structure that holds files and (smaller) directories”, or “a family tree starts with a couple who have children, each with their own family sub-trees”.

Programming languages generally support **recursion**, which means that, in order to solve a problem, functions can *call themselves* to solve smaller subproblems.

Any problem that can be solved iteratively (with a for or while loop) can also be solved recursively. However, recursion takes a while wrap your head around, and because of this, it is generally only used in specific cases, where either your problem is recursive in nature, or your data is recursive.

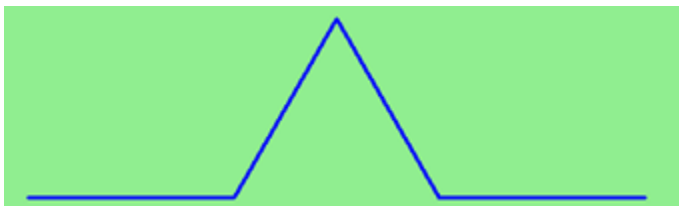
### 10.1 Drawing Fractals

For our purposes, a **fractal** is a drawing which also has *self-similar* structure, where it can be defined in terms of itself. This is a typical example of a problem which is recursive in nature.

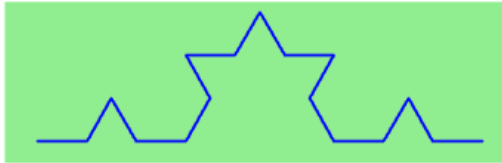
Let us start by looking at the famous Koch fractal. An order 0 Koch fractal is simply a straight line of a given size.



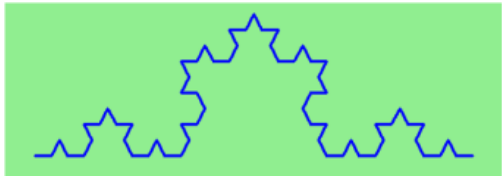
An order 1 Koch fractal is obtained like this: instead of drawing just one line, draw instead four smaller segments, in the pattern shown here:



Now what would happen if we repeated this Koch pattern again on each of the order 1 segments? We'd get this order 2 Koch fractal:



Repeating our pattern again gets us an order 3 Koch fractal:



Now let us think about it the other way around. To draw a Koch fractal of order 3, we can simply draw four order 2 Koch fractals. But each of these in turn needs four order 1 Koch fractals, and each of those in turn needs four order 0 fractals. Ultimately, the only drawing that will take place is at order 0. This is very simple to code up in Python:

```
1 def koch(tortoise, order, size):
2     """
3     Make turtle tortoise draw a Koch fractal of 'order' and 'size'.
4     Leave the turtle facing the same direction.
5     """
6
7     if order == 0:          # The base case is just a straight line
8         tortoise.forward(size)
9     else:
10        koch(tortoise, order-1, size/3)    # Go 1/3 of the way
11        tortoise.left(60)
12        koch(tortoise, order-1, size/3)
13        tortoise.right(120)
14        koch(tortoise, order-1, size/3)
15        tortoise.left(60)
16        koch(tortoise, order-1, size/3)
```

The key thing that is new here is that if order is not zero, koch calls itself recursively to get its job done.

Let's make a simple observation and tighten up this code. Remember that turning right by 120 is the same as turning left by -120. So with a bit of clever rearrangement, we can use a loop instead of lines 10-16:

```
1 def koch(tortoise, order, size):
2     if order == 0:
3         tortoise.forward(size)
4     else:
5         for angle in [60, -120, 60, 0]:
6             koch(tortoise, order-1, size/3)
7             tortoise.left(angle)
```

The final turn is 0 degrees — so it has no effect. But it has allowed us to find a pattern and reduce seven lines of code to three, which will make things easier for our next observations.

---

### Recursion, the high-level view

One way to think about this is to convince yourself that the function works correctly when you call it for an order 0 fractal. Then do a mental *leap of faith*, saying “*the fairy godmother*” (or Python, if you can think of Python as your

fairy godmother) *knows how to handle the recursive level 0 calls for me on lines 11, 13, 15, and 17, so I don't need to think about that detail!*" All I need to focus on is how to draw an order 1 fractal *if I can assume the order 0 one is already working*.

You're practicing *mental abstraction* — ignoring the subproblem while you solve the big problem.

If this mode of thinking works (and you should practice it!), then take it to the next level. Aha! now can I see that it will work when called for order 2 *under the assumption that it is already working for level 1*.

And, in general, if I can assume the order  $n-1$  case works, can I just solve the level  $n$  problem?

Students of mathematics who have played with proofs of induction should see some very strong similarities here.

---

### Recursion, the low-level operational view

Another way of trying to understand recursion is to get rid of it! If we had separate functions to draw a level 3 fractal, a level 2 fractal, a level 1 fractal and a level 0 fractal, we could simplify the above code, quite mechanically, to a situation where there was no longer any recursion, like this:

```
1 def koch_0(tortoise, size):
2     tortoise.forward(size)
3
4 def koch_1(tortoise, size):
5     for angle in [60, -120, 60, 0]:
6         koch_0(tortoise, size/3)
7         tortoise.left(angle)
8
9 def koch_2(tortoise, size):
10    for angle in [60, -120, 60, 0]:
11        koch_1(tortoise, size/3)
12        tortoise.left(angle)
13
14 def koch_3(tortoise, size):
15    for angle in [60, -120, 60, 0]:
16        koch_2(tortoise, size/3)
17        tortoise.left(angle)
```

This trick of “unrolling” the recursion gives us an operational view of what happens. You can trace the program into `koch_3`, and from there, into `koch_2`, and then into `koch_1`, etc., all the way down the different layers of the recursion.

This might be a useful hint to build your understanding. The mental goal is, however, to be able to do the abstraction!

---

## 10.2 Recursive data structures

Most of the Python data types we have seen can be grouped inside lists and tuples in a variety of ways. Lists and tuples can also be nested, providing many possibilities for organizing data. The organization of data for the purpose of making it easier to use is called a **data structure**.

It's election time and we are helping to compute the votes as they come in. Votes arriving from individual wards, precincts, municipalities, counties, and states are sometimes reported as a sum total of votes and sometimes as a list of subtotals of votes. After considering how best to store the tallies, we decide to use a *nested number list*, which we define as follows:

A *nested number list* is a list whose elements are either:

1. numbers
2. nested number lists

Notice that the term, *nested number list* is used in its own definition. **Recursive definitions** like this are quite common in mathematics and computer science. They provide a concise and powerful way to describe **recursive data structures** that are partially composed of smaller and simpler instances of themselves. The definition is not circular, since at some point we will reach a list that does not have any lists as elements.

Now suppose our job is to write a function that will sum all of the values in a nested number list. Python has a built-in function which finds the sum of a sequence of numbers:

```
>>> sum([1, 2, 8])
11
```

For our *nested number list*, however, `sum` will not work:

```
>>> sum([1, 2, [11, 13], 8])
Traceback (most recent call last):
  File "<interactive input>", line 1, in <module>
TypeError: unsupported operand type(s) for +: 'int' and 'list'
>>>
```

The problem is that the third element of this list, `[11, 13]`, is itself a list, so it cannot just be added to 1, 2, and 8.

## 10.3 Processing recursive number lists

To sum all the numbers in our recursive nested number list we need to traverse the list, visiting each of the elements within its nested structure, adding any numeric elements to our sum, and *recursively repeating the summing process* with any elements which are themselves sub-lists.

Thanks to recursion, the Python code needed to sum the values of a nested number list is surprisingly short:

```
1 def recursive_sum(nested_number_list):
2     """Returns the total sum of all elements in nested_number_list"""
3     total = 0
4     for element in nested_number_list:
5         if type(element) is list:
6             total += recursive_sum(element)
7         else:
8             total += element
9     return total
```

The body of `recursive_sum` consists mainly of a `for` loop that traverses `nested_number_list`. If `element` is a numerical value (the `else` branch), it is simply added to `total`. If `element` is a list, then `recursive_sum` is called again, with the `element` as an argument. The statement inside the function definition in which the function calls itself is known as the **recursive call**.

The example above has a **base case** (on line 13) which does not lead to a recursive call: the case where the `element` is not a (sub-) list. Without a base case, you'll have **infinite recursion**, and your program will not work.

An alternative solution, completely recursive, would be the following. Notice that this implementation does not contain a `for` loop!

```
1 def recursive_sum(nested_number_list):
2     """Returns the total sum of all elements in nested_number_list"""
3     if len(nested_number_list) == 0:
```

(continues on next page)

(continued from previous page)

```
4     return 0
5     head, *tail = nested_number_list #Assign the first element of nested_number_list
    ↳to head, and the rest to tail.
6     if isinstance(head, list): # If head is a list....
7         return recursive_sum(head) + recursive_sum(tail)
8     else:
9         return head + recursive_sum(tail)
```

Recursion is truly one of the most beautiful and elegant tools in computer science.

A slightly more complicated problem is finding the largest value in our nested number list:

```
1 def recursive_max(nested_list):
2     """
3     Find the maximum in a recursive structure of lists
4     within other lists.
5     Precondition: No lists or sublists are empty.
6     """
7     largest = None
8     first_time = True
9     for element in nested_list:
10        if type(element) is list:
11            value = recursive_max(element)
12        else:
13            value = element
14
15        if first_time or value > largest:
16            largest = value
17            first_time = False
18
19    return largest
```

The added twist to this problem is finding a value for initializing `largest`. We can't just use `nested_list[0]`, since that could be either a element or a list. To solve this problem (at every recursive call) we initialize a Boolean flag (at line 8). When we've found the value of interest, (at line 15) we check to see whether this is the initializing (first) value for `largest`, or a value that could potentially change `largest`.

Again here we have a base case at line 13. If we don't supply a base case, Python stops after reaching a maximum recursion depth and returns a runtime error. See how this happens, by running this little script which we will call *infinite\_recursion.py*:

```
1 def recursion_depth(number):
2     print("{0}, ".format(number), end="")
3     recursion_depth(number + 1)
4
5 recursion_depth(0)
```

After watching the messages flash by, you will be presented with the end of a long traceback that ends with a message like the following:

```
RuntimeError: maximum recursion depth exceeded ...
```

We would certainly never want something like this to happen to a user of one of our programs, so in another appendix we'll see how errors, any kinds of errors, are handled in Python.



## 10.4 Case study: Fibonacci numbers

The famous **Fibonacci sequence** 0, 1, 1, 2, 3, 5, 8, 13, 21, 34, 55, 89, 134, ... was devised by Fibonacci (1170-1250), who used this to model the breeding of (pairs) of rabbits. If, in generation 7 you had 21 pairs in total, of which 13 were adults, then next generation the adults will all have bred new children, and the previous children will have grown up to become adults. So in generation 8 you'll have 13+21=34, of which 21 are adults.

This *model* to explain rabbit breeding made the simplifying assumption that rabbits never died. Scientists often make (unrealistic) simplifying assumptions and restrictions to make some headway with the problem.

If we number the terms of the sequence from 0, we can describe each term recursively as the sum of the previous two terms:

```
fib(0) = 0
fib(1) = 1
fib(n) = fib(n-1) + fib(n-2)   for n >= 2
```

This translates very directly into some Python:

```
1 def fib(n):
2     if n <= 1:
3         return n
4     t = fib(n-1) + fib(n-2)
5     return t
```

This is a particularly inefficient algorithm, and this could be solved far more efficient iteratively:

```
1 import time
2 t0 = time.clock()
3 n = 35
4 result = fib(n)
5 t1 = time.clock()
6
7 print("fib({0}) = {1}, ({2:.2f} secs)".format(n, result, t1-t0))
```

We get the correct result, but an exploding amount of work!

```
fib(35) = 9227465, (10.54 secs)
```

## 10.5 Example with recursive directories and files

The following program lists the contents of a directory and all its subdirectories.

```
1 import os
2
3 def get_dirlist(path):
4     """
5     Return a sorted list of all entries in path.
6     This returns just the names, not the full path to the names.
7     """
8     dirlist = os.listdir(path)
9     dirlist.sort()
10    return dirlist
11
12 def print_files(path, prefix = ""):
```

(continues on next page)

(continued from previous page)

```
13 """ Print recursive listing of contents of path """
14 if prefix == "": # Detect outermost call, print a heading
15     print("Folder listing for", path)
16     prefix = "| "
17
18     dirlist = get_dirlist(path)
19     for file in dirlist:
20         print(prefix+file)          # Print the line
21         fullname = os.path.join(path, file) # Turn name into full pathname
22         if os.path.isdir(fullname):      # If a directory, recurse.
23             print_files(fullname, prefix + "| ")
```

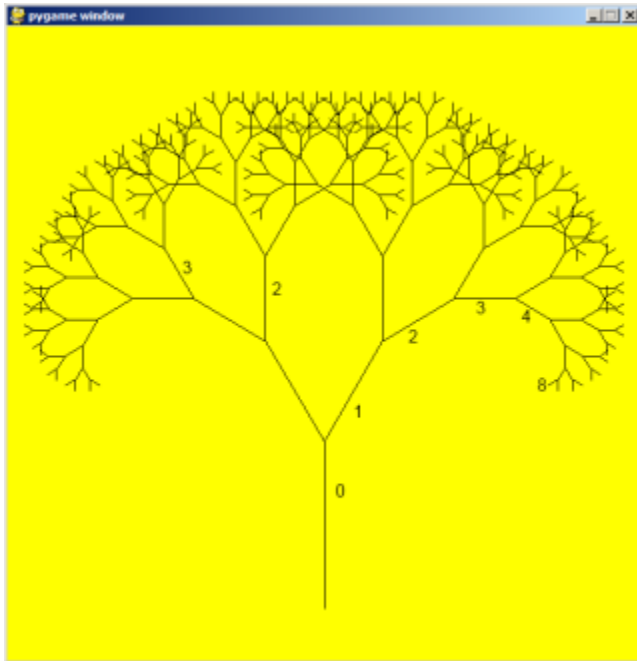
Calling the function `print_files` with some folder name will produce output similar to this:

```
Folder listing for c:\python31\Lib\site-packages\pygame\examples
| __init__.py
| aacircle.py
| aliens.py
| arraydemo.py
| blend_fill.py
| blit_blends.py
| camera.py
| chimp.py
| cursors.py
| data
| | alien1.png
| | alien2.png
| | alien3.png
...
```

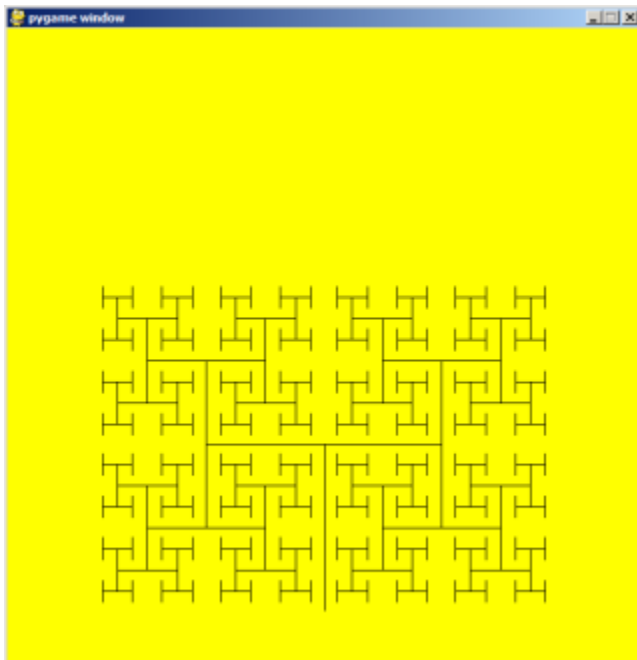
Note that something similar is already implemented in the `os` module: `os.walk`.

## 10.6 An animated fractal, using PyGame

Here we have a tree fractal pattern of order 8. We've labelled some of the edges, showing the depth of the recursion at which each edge was drawn.



In the tree above, the angle of deviation from the trunk is 30 degrees. Varying that angle gives other interesting shapes, for example, with the angle at 90 degrees we get this:



An interesting animation occurs if we generate and draw trees very rapidly, each time varying the angle a little. Although the Turtle module can draw trees like this quite elegantly, we could struggle for good frame rates. So we'll use PyGame instead, with a few embellishments and observations. (Once again, we suggest you cut and paste this code into your Python environment.)

```
1 import pygame, math
2 pygame.init()           # prepare the pygame module for use
3
```

(continues on next page)

(continued from previous page)

```

4  # Create a new surface and window.
5  surface_size = 1024
6  main_surface = pygame.display.set_mode((surface_size,surface_size))
7  my_clock = pygame.time.Clock()
8
9
10 def draw_tree(order, theta, size, position, heading, color=(0,0,0), depth=0):
11
12     trunk_ratio = 0.29      # How big is the trunk relative to whole tree?
13     trunk = size * trunk_ratio # length of trunk
14     delta_x = trunk * math.cos(heading)
15     delta_y = trunk * math.sin(heading)
16     (u, v) = position
17     newposition = (u + delta_x, v + delta_y)
18     pygame.draw.line(main_surface, color, position, newposition)
19
20     if order > 0:    # Draw another layer of subtrees
21
22         # These next six lines are a simple hack to make the two major halves
23         # of the recursion different colors. Fiddle here to change colors
24         # at other depths, or when depth is even, or odd, etc.
25         if depth == 0:
26             color1 = (255, 0, 0)
27             color2 = (0, 0, 255)
28         else:
29             color1 = color
30             color2 = color
31
32         # make the recursive calls to draw the two subtrees
33         newsize = size*(1 - trunk_ratio)
34         draw_tree(order-1, theta, newsize, newposition, heading-theta, color1, depth+1)
35         draw_tree(order-1, theta, newsize, newposition, heading+theta, color2, depth+1)
36
37
38 def gameloop():
39
40     theta = 0
41     while True:
42
43         # Handle events from keyboard, mouse, etc.
44         event = pygame.event.poll()
45         if event.type == pygame.QUIT:
46             break;
47
48         # Updates - change the angle
49         theta += 0.01
50
51         # Draw everything
52         main_surface.fill((255, 255, 0))
53         draw_tree(9, theta, surface_size*0.9, (surface_size//2, surface_size-50), -
↳math.pi/2)
54
55         pygame.display.flip()
56         my_clock.tick(120)
57
58
59 gameloop()

```

(continues on next page)

(continued from previous page)

```
60 pygame.quit()
```

- The `math` library works with angles in radians rather than degrees.
- Lines 14 and 15 uses some high school trigonometry. From the length of the desired line (`trunk`), and its desired angle, `cos` and `sin` help us to calculate the `x` and `y` distances we need to move.
- Lines 22-30 are unnecessary, except if we want a colorful tree.
- In the main game loop at line 49 we change the angle on every frame, and redraw the new tree.
- Line 18 shows that PyGame can also draw lines, and plenty more. Check out the documentation. For example, drawing a small circle at each branch point of the tree can be accomplished by adding this line directly below line 18:

```
1 pygame.draw.circle(main_surface, color, (int(position[0]), int(position[1])), 3)
```

Another interesting effect — instructive too, if you wish to reinforce the idea of different instances of the function being called at different depths of recursion — is to create a list of colors, and let each recursive depth use a different color for drawing. (Use the depth of the recursion to index the list of colors.)

## 10.7 Mutual Recursion

In addition to a function calling just itself, it is also possible to make multiple functions that call eachother. This is rarely really usefull, but it can be used to make state machines.

```
1 def function_a(n): # Do things associated with state A
2     if n == 0:
3         return
4     print('a')
5     function_b(n - 1) # Proceed to state B
6
7
8 def function_b(n): # Do things associated with state B
9     print('b')
10    function_a(n - 1) # Proceed to state A
```

## 10.8 Glossary

**base case** A branch of the conditional statement in a recursive function that does not give rise to further recursive calls.

**infinite recursion** A function that calls itself recursively without ever reaching any base case. Eventually, infinite recursion causes a runtime error.

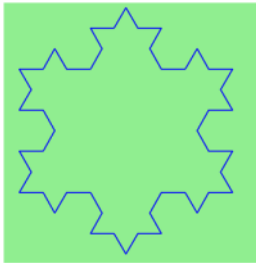
**recursion** The process of calling a function that is already executing.

**recursive call** The statement that calls an already executing function. Recursion can also be indirect — function *f* can call *g* which calls *h*, and *h* could make a call back to *f*.

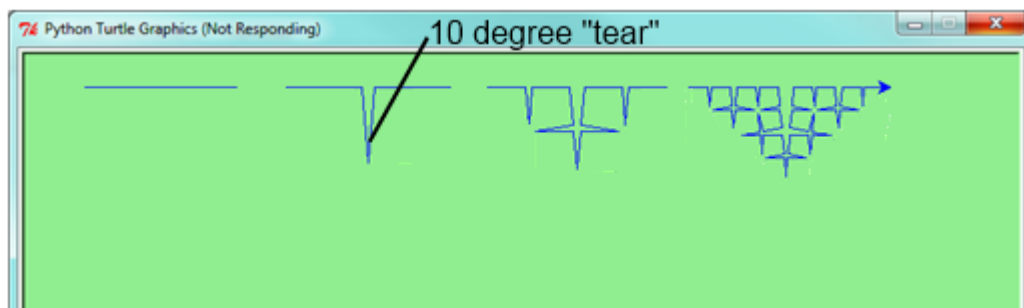
**recursive definition** A definition which defines something in terms of itself. To be useful it must include *base cases* which are not recursive. In this way it differs from a *circular definition*. Recursive definitions often provide an elegant way to express complex data structures, like a directory that can contain other directories, or a menu that can contain other menus.

## 10.9 Exercises

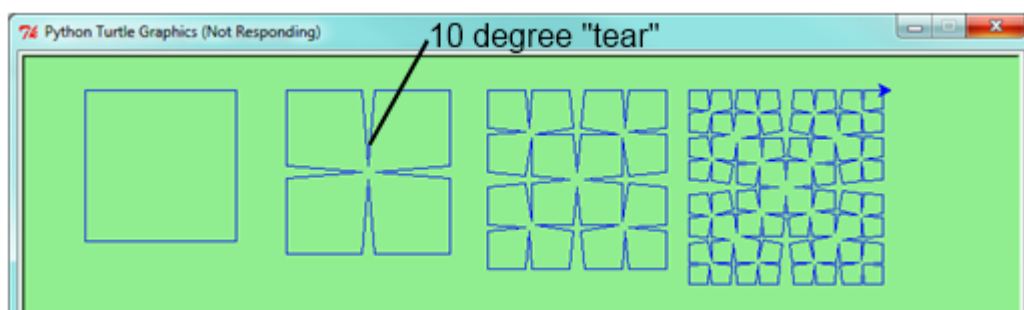
1. Modify the Koch fractal program so that it draws a Koch snowflake, like this:



2. (a) Draw a Cesaro torn line fractal, of the order given by the user. We show four different lines of orders 0,1,2,3. In this example, the angle of the tear is 10 degrees.



- (b) Four lines make a square. Use the code in part a) to draw cesaro squares. Varying the angle gives interesting effects — experiment a bit, or perhaps let the user input the angle of the tear.

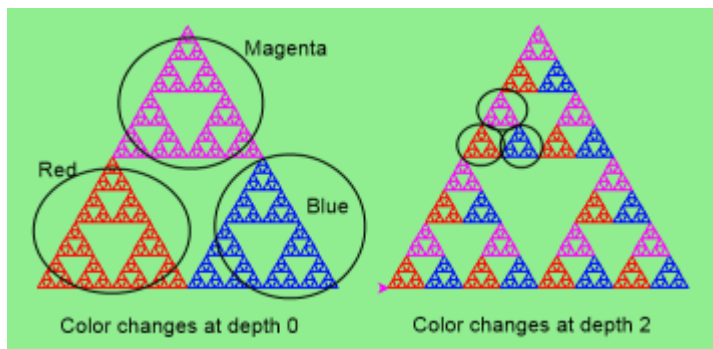


- (c) (For the mathematically inclined). In the squares shown here, the higher-order drawings become a little larger. (Look at the bottom lines of each square - they're not aligned.) This is because we just halved the drawn part of the line for each recursive subproblem. So we've "grown" the overall square by the width of the tear(s). Can you solve the geometry problem so that the total size of the subproblem case (including the tear) remains exactly the same size as the original?

3. A Sierpinski triangle of order 0 is an equilateral triangle. An order 1 triangle can be drawn by drawing 3 smaller triangles (shown slightly disconnected here, just to help our understanding). Higher order 2 and 3 triangles are also shown. Draw Sierpinski triangles of any order input by the user.



4. Adapt the above program to change the color of its three sub-triangles at some depth of recursion. The illustration below shows two cases: on the left, the color is changed at depth 0 (the outmost level of recursion), on the right, at depth 2. If the user supplies a negative depth, the color never changes. (Hint: add a new optional parameter `colorChangeDepth` (which defaults to -1), and make this one smaller on each recursive subcall. Then, in the section of code before you recurse, test whether the parameter is zero, and change color.)



5. Write a function, `recursive_min`, that returns the smallest value in a nested number list. Assume there are no empty lists or sublists:
6. Write a function `count` that returns the number of occurrences of `target` in a nested list:
7. Write a function `flatten` that returns a simple list containing all the values in a nested list:
8. Rewrite the fibonacci algorithm without using recursion. Can you find bigger terms of the sequence? Can you find `fib(200)`?
9. Use help to find out what `sys.getrecursionlimit()` and `sys.setrecursionlimit(n)` do. Create several experiments similar to what was done in *infinite\_recursion.py* to test your understanding of how these module functions work.
10. Write a program that walks a directory structure (as in the last section of this chapter), but instead of printing filenames, it returns a list of all the full paths of files in the directory or the subdirectories. (Don't include directories in this list — just files.) For example, the output list might have elements like this:

```
[ "C:\\Python31\\Lib\\site-packages\\pygame\\docs\\ref\\mask.html",  
  "C:\\Python31\\Lib\\site-packages\\pygame\\docs\\ref\\midi.html",  
  ...  
  "C:\\Python31\\Lib\\site-packages\\pygame\\examples\\aliens.py",  
  ...  
  "C:\\Python31\\Lib\\site-packages\\pygame\\examples\\data\\boom.wav",  
  ... ]
```

11. Write a program named `litter.py` that creates an empty file named `trash.txt` in each subdirectory of a directory tree given the root of the tree as an argument (or the current directory as a default). Now write a program named `cleanup.py` that removes all these files.

*Hint #1:* Use the program from the example in the last section of this chapter as a basis for these two recursive programs. Because you're going to destroy files on your disks, you better get this right, or you risk losing files you care about. So excellent advice is that initially you should fake the deletion of the files — just print the full path names of each file that you intend to delete. Once you're happy that your logic is correct, and you can see that you're not deleting the wrong things, you can replace the print statement with the real thing.

*Hint #2:* Look in the `os` module for a function that removes files.





### 11.1 Classes and Objects — the Basics

#### 11.1.1 Object-oriented programming

Python is an **object-oriented programming language**, which means that it provides features that support **object-oriented programming (OOP)**.

Object-oriented programming has its roots in the 1960s, but it wasn't until the mid 1980s that it became the main **programming paradigm** used in the creation of new software. It was developed as a way to handle the rapidly increasing size and complexity of software systems, and to make it easier to modify these large and complex systems over time.

Up to now, most of the programs we have been writing use a **procedural programming** paradigm. In procedural programming the focus is on writing functions or *procedures* which operate on data. In object-oriented programming the focus is on the creation of **objects** which contain both data and functionality together. (We have seen turtle objects, string objects, and random number generators, to name a few places where we've already worked with objects.)

Usually, each object definition corresponds to some object or concept in the real world, and the functions that operate on that object correspond to the ways real-world objects interact.

#### 11.1.2 User-defined compound data types

We've already seen classes like `str`, `int`, `float` and `Turtle`. We are now ready to create our own user-defined class: the `Point`.

Consider the concept of a mathematical point. In two dimensions, a point is two numbers (coordinates) that are treated collectively as a single object. Points are often written in between parentheses with a comma separating the coordinates. For example, `(0, 0)` represents the origin, and `(x, y)` represents the point `x` units to the right and `y` units up from the origin.

Some of the typical operations that one associates with points might be calculating the distance of a point from the origin, or from another point, or finding a midpoint of two points, or asking if a point falls within a given rectangle or circle. We'll shortly see how we can organize these together with the data.

A natural way to represent a point in Python is with two numeric values. The question, then, is how to group these two values into a compound object. The quick and dirty solution is to use a tuple, and for some applications that might be a good choice.

An alternative is to define a new **class**. This approach involves a bit more effort, but it has advantages that will be apparent soon. We'll want our points to each have an *x* and a *y* attribute, so our first class definition looks like this:

```
1 class Point:
2     """ Point class represents and manipulates x,y coords. """
3
4     def __init__(self):
5         """ Create a new point at the origin """
6         self.x = 0
7         self.y = 0
```

Class definitions can appear anywhere in a program, but they are usually near the beginning (after the `import` statements). Some programmers and languages prefer to put every class in a module of its own — we won't do that here. The syntax rules for a class definition are the same as for other compound statements. There is a header which begins with the keyword, `class`, followed by the name of the class, and ending with a colon. Indentation levels tell us where the class ends.

If the first line after the class header is a string, it becomes the docstring of the class, and will be recognized by various tools. (This is also the way docstrings work in functions.)

Every class should have a method with the special name `__init__`. This **initializer method** is automatically called whenever a new instance of `Point` is created. It gives the programmer the opportunity to set up the attributes required within the new instance by giving them their initial state/values. The `self` parameter (we could choose any other name, but `self` is the convention) is automatically set to reference the newly created object that needs to be initialized.

So let's use our new `Point` class now:

```
1 p = Point()          # Instantiate an object of type Point
2 q = Point()          # Make a second point
3
4 print(p.x, p.y, q.x, q.y) # Each point object has its own x and y
```

This program prints:

```
0 0 0 0
```

because during the initialization of the objects, we created two attributes called *x* and *y* for each, and gave them both the value 0.

This should look familiar — we've used classes before to create more than one object:

```
1 from turtle import Turtle
2
3 tess = Turtle()      # Instantiate objects of type Turtle
4 alex = Turtle()
```

The variables `p` and `q` are assigned references to two new `Point` objects. A function like `Turtle` or `Point` that creates a new object instance is called a **constructor**, and every class automatically provides a constructor function which is named the same as the class.

It may be helpful to think of a class as a *factory* for making objects. The class itself isn't an instance of a point, but it contains the machinery to make point instances. Every time we call the constructor, we're asking the factory to make us a new object. As the object comes off the production line, its initialization method is executed to get the object properly set up with its factory default settings.

The combined process of “make me a new object” and “get its settings initialized to the factory default settings” is called **instantiation**.

### 11.1.3 Attributes

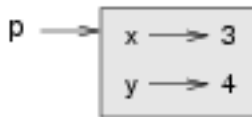
Like real world objects, object instances have both attributes and methods.

We can modify the attributes in an instance using dot notation:

```
>>> p.x = 3
>>> p.y = 4
```

Both modules and instances create their own namespaces, and the syntax for accessing names contained in each, called **attributes**, is the same. In this case the attribute we are selecting is a data item from an instance.

The following state diagram shows the result of these assignments:



The variable `p` refers to a `Point` object, which contains two attributes. Each attribute refers to a number.

We can access the value of an attribute using the same syntax:

```
>>> print(p.y)
4
>>> x = p.x
>>> print(x)
3
```

The expression `p.x` means, “Go to the object `p` refers to and get the value of `x`”. In this case, we assign that value to a variable named `x`. There is no conflict between the variable `x` (in the global namespace here) and the attribute `x` (in the namespace belonging to the instance). The purpose of dot notation is to fully qualify which variable we are referring to unambiguously.

We can use dot notation as part of any expression, so the following statements are legal:

```
1 print("(x={0}, y={1})".format(p.x, p.y))
2 distance_squared_from_origin = p.x * p.x + p.y * p.y
```

The first line outputs `(x=3, y=4)`. The second line calculates the value 25.

### 11.1.4 Improving our initializer

To create a point at position (7, 6) currently needs three lines of code:

```
1 p = Point()
2 p.x = 7
3 p.y = 6
```

We can make our class constructor more general by placing extra parameters into the `__init__` method, as shown in this example:

```
1 class Point:
2     """ Point class represents and manipulates x,y coords. """
3
4     def __init__(self, x=0, y=0):
5         """ Create a new point at x, y """
6         self.x = x
7         self.y = y
8
9     # Other statements outside the class continue below here.
```

The `x` and `y` parameters here are both optional. If the caller does not supply arguments, they'll get the default values of 0. Here is our improved class in action:

```
>>> p = Point(4, 2)
>>> q = Point(6, 3)
>>> r = Point()           # r represents the origin (0, 0)
>>> print(p.x, q.y, r.x)
4 3 0
```

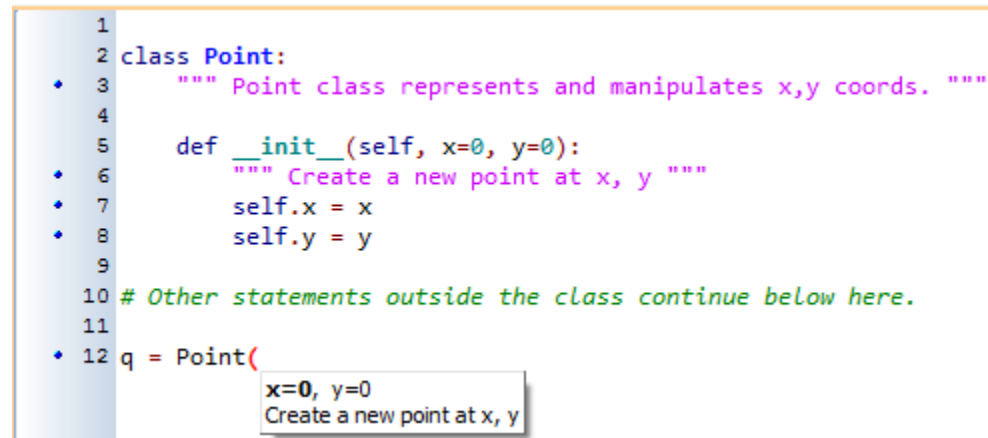
---

### Technically speaking ...

If we are really fussy, we would argue that the `__init__` method's docstring is inaccurate. `__init__` doesn't *create* the object (i.e. set aside memory for it), — it just initializes the object to its factory-default settings after its creation.

But tools like PyScripter understand that instantiation — creation and initialization — happen together, and they choose to display the *initializer's* docstring as the tooltip to guide the programmer that calls the class constructor.

So we're writing the docstring so that it makes the most sense when it pops up to help the programmer who is using our `Point` class:



```
1
2 class Point:
3     """ Point class represents and manipulates x,y coords. """
4
5     def __init__(self, x=0, y=0):
6         """ Create a new point at x, y """
7         self.x = x
8         self.y = y
9
10    # Other statements outside the class continue below here.
11
12    q = Point(
```

x=0, y=0  
Create a new point at x, y

### 11.1.5 Adding other methods to our class

The key advantage of using a class like `Point` rather than a simple tuple `(6, 7)` now becomes apparent. We can add methods to the `Point` class that are sensible operations for points, but which may not be appropriate for other tuples like `(25, 12)` which might represent, say, a day and a month, e.g. Christmas day. So being able to calculate the distance from the origin is sensible for points, but not for (day, month) data. For (day, month) data, we'd like different operations, perhaps to find what day of the week it will fall on in 2020.

Creating a class like `Point` brings an exceptional amount of “organizational power” to our programs, and to our thinking. We can group together the sensible operations, and the kinds of data they apply to, and each instance of the class can have its own state.

A **method** behaves like a function but it is invoked on a specific instance, e.g. `tess.right(90)`. Like a data attribute, methods are accessed using dot notation.

Let’s add another method, `distance_from_origin`, to see better how methods work:

```
1 class Point:
2     """ Create a new Point, at coordinates x, y """
3
4     def __init__(self, x=0, y=0):
5         """ Create a new point at x, y """
6         self.x = x
7         self.y = y
8
9     def distance_from_origin(self):
10        """ Compute my distance from the origin """
11        return ((self.x ** 2) + (self.y ** 2)) ** 0.5
```

Let’s create a few point instances, look at their attributes, and call our new method on them: (We must run our program first, to make our `Point` class available to the interpreter.)

```
>>> p = Point(3, 4)
>>> p.x
3
>>> p.y
4
>>> p.distance_from_origin()
5.0
>>> q = Point(5, 12)
>>> q.x
5
>>> q.y
12
>>> q.distance_from_origin()
13.0
>>> r = Point()
>>> r.x
0
>>> r.y
0
>>> r.distance_from_origin()
0.0
```

When defining a method, the first parameter refers to the instance being manipulated. As already noted, it is customary to name this parameter `self`.

Notice that the caller of `distance_from_origin` does not explicitly supply an argument to match the `self` parameter — this is done for us, behind our back.

### 11.1.6 Instances as arguments and parameters

We can pass an object as an argument in the usual way. We’ve already seen this in some of the turtle examples, where we passed the turtle to some function like `draw_bar` in the chapter titled *Conditionals*, so that the function could control and use whatever turtle instance we passed to it.

Be aware that our variable only holds a reference to an object, so passing `tess` into a function creates an alias: both the caller and the called function now have a reference, but there is only one turtle!

Here is a simple function involving our new `Point` objects:

```
1 def print_point(pt):
2     print("{0}, {1}".format(pt.x, pt.y))
```

`print_point` takes a point as an argument and formats the output in whichever way we choose. If we call `print_point(p)` with point `p` as defined previously, the output is `(3, 4)`.

### 11.1.7 Converting an instance to a string

Most object-oriented programmers probably would not do what we've just done in `print_point`. When we're working with classes and objects, a preferred alternative is to add a new method to the class. And we don't like chatterbox methods that call `print`. A better approach is to have a method so that every instance can produce a string representation of itself. Let's initially call it `to_string`:

```
1 class Point:
2     # ...
3
4     def to_string(self):
5         return "{0}, {1}".format(self.x, self.y)
```

Now we can say:

```
>>> p = Point(3, 4)
>>> print(p.to_string())
(3, 4)
```

But don't we already have a `str` type converter that can turn our object into a string? Yes! And doesn't `print` automatically use this when printing things? Yes again! But these automatic mechanisms do not yet do exactly what we want:

```
>>> str(p)
'<__main__.Point object at 0x01F9AA10>'
>>> print(p)
'<__main__.Point object at 0x01F9AA10>'
```

Python has a clever trick up its sleeve to fix this. If we call our new method `__str__` instead of `to_string`, the Python interpreter will use our code whenever it needs to convert a `Point` to a string. Let's re-do this again, now:

```
1 class Point:
2     # ...
3
4     def __str__(self):      # All we have done is renamed the method
5         return "{0}, {1}".format(self.x, self.y)
```

and now things are looking great!

```
>>> str(p)      # Python now uses the __str__ method that we wrote.
(3, 4)
>>> print(p)
(3, 4)
```

### 11.1.8 Instances as return values

Functions and methods can return instances. For example, given two `Point` objects, find their midpoint. First we'll write this as a regular function:

```
1 def midpoint(p1, p2):
2     """ Return the midpoint of points p1 and p2 """
3     mx = (p1.x + p2.x)/2
4     my = (p1.y + p2.y)/2
5     return Point(mx, my)
```

The function creates and returns a new `Point` object:

```
>>> p = Point(3, 4)
>>> q = Point(5, 12)
>>> r = midpoint(p, q)
>>> r
(4.0, 8.0)
```

Now let us do this as a method instead. Suppose we have a point object, and wish to find the midpoint halfway between it and some other target point:

```
1 class Point:
2     # ...
3
4     def halfway(self, target):
5         """ Return the halfway point between myself and the target """
6         mx = (self.x + target.x)/2
7         my = (self.y + target.y)/2
8         return Point(mx, my)
```

This method is identical to the function, aside from some renaming. It's usage might be like this:

```
>>> p = Point(3, 4)
>>> q = Point(5, 12)
>>> r = p.halfway(q)
>>> r
(4.0, 8.0)
```

While this example assigns each point to a variable, this need not be done. Just as function calls are composable, method calls and object instantiation are also composable, leading to this alternative that uses no variables:

```
>>> print(Point(3, 4).halfway(Point(5, 12)))
(4.0, 8.0)
```

### 11.1.9 A change of perspective

The original syntax for a function call, `print_time(current_time)`, suggests that the function is the active agent. It says something like, “Hey, `print_time`! Here’s an object for you to print.”

In object-oriented programming, the objects are considered the active agents. An invocation like `current_time.print_time()` says “Hey `current_time`! Please print yourself!”

In our early introduction to turtles, we used an object-oriented style, so that we said `tess.forward(100)`, which asks the turtle to move itself forward by the given number of steps.



This change in perspective might be more polite, but it may not initially be obvious that it is useful. But sometimes shifting responsibility from the functions onto the objects makes it possible to write more versatile functions, and makes it easier to maintain and reuse code.

The most important advantage of the object-oriented style is that it fits our mental chunking and real-life experience more accurately. In real life our `cook` method is part of our microwave oven — we don't have a `cook` function sitting in the corner of the kitchen, into which we pass the microwave! Similarly, we use the cellphone's own methods to send an sms, or to change its state to silent. The functionality of real-world objects tends to be tightly bound up inside the objects themselves. OOP allows us to accurately mirror this when we organize our programs.

### 11.1.10 Objects can have state

Objects are most useful when we also need to keep some state that is updated from time to time. Consider a turtle object. Its state consists of things like its position, its heading, its color, and its shape. A method like `left(90)` updates the turtle's heading, `forward` changes its position, and so on.

For a bank account object, a main component of the state would be the current balance, and perhaps a log of all transactions. The methods would allow us to query the current balance, deposit new funds, or make a payment. Making a payment would include an amount, and a description, so that this could be added to the transaction log. We'd also want a method to show the transaction log.

### 11.1.11 Glossary

**attribute** One of the named data items that makes up an instance.

**class** A user-defined compound type. A class can also be thought of as a template for the objects that are instances of it. (The iPhone is a class. By December 2010, estimates are that 50 million instances had been sold!)

**constructor** Every class has a “factory”, called by the same name as the class, for making new instances. If the class has an *initializer method*, this method is used to get the attributes (i.e. the state) of the new object properly set up.

**initializer method** A special method in Python (called `__init__`) that is invoked automatically to set a newly created object's attributes to their initial (factory-default) state.

**instance** An object whose type is of some class. Instance and object are used interchangeably.

**instantiate** To create an instance of a class, and to run its initializer.

**method** A function that is defined inside a class definition and is invoked on instances of that class.

**object** A compound data type that is often used to model a thing or concept in the real world. It bundles together the data and the operations that are relevant for that kind of data. Instance and object are used interchangeably.

**object-oriented programming** A powerful style of programming in which data and the operations that manipulate it are organized into objects.

**object-oriented language** A language that provides features, such as user-defined classes and inheritance, that facilitate object-oriented programming.

### 11.1.12 Exercises

1. Rewrite the `distance` function from the chapter titled *Fruitful functions* so that it takes two `Points` as parameters instead of four numbers.
2. Add a method `reflect_x` to `Point` which returns a new `Point`, one which is the reflection of the point about the x-axis. For example, `Point(3, 5).reflect_x()` is `(3, -5)`

3. Add a method `slope_from_origin` which returns the slope of the line joining the origin to the point. For example,

```
>>> Point(4, 10).slope_from_origin()
2.5
```

What cases will cause this method to fail?

4. The equation of a straight line is “ $y = ax + b$ ”, (or perhaps “ $y = mx + c$ ”). The coefficients  $a$  and  $b$  completely describe the line. Write a method in the `Point` class so that if a point instance is given another point, it will compute the equation of the straight line joining the two points. It must return the two coefficients as a tuple of two values. For example,

```
>>> print(Point(4, 11).get_line_to(Point(6, 15)))
>>> (2, 3)
```

This tells us that the equation of the line joining the two points is “ $y = 2x + 3$ ”. When will this method fail?

5. Given four points that fall on the circumference of a circle, find the midpoint of the circle. When will this function fail?

*Hint:* You *must* know how to solve the geometry problem *before* you think of going anywhere near programming. You cannot program a solution to a problem if you don’t understand what you want the computer to do!

6. Create a new class, `SMS_store`. The class will instantiate `SMS_store` objects, similar to an inbox or outbox on a cellphone:

```
my_inbox = SMS_store()
```

This store can hold multiple SMS messages (i.e. its internal state will just be a list of messages). Each message will be represented as a tuple:

```
(has_been_viewed, from_number, time_arrived, text_of_SMS)
```

The inbox object should provide these methods:

```
my_inbox.add_new_arrival(from_number, time_arrived, text_of_SMS)
    # Makes new SMS tuple, inserts it after other messages
    # in the store. When creating this message, its
    # has_been_viewed status is set False.

my_inbox.message_count()
    # Returns the number of sms messages in my_inbox

my_inbox.get_unread_indexes()
    # Returns list of indexes of all not-yet-viewed SMS messages

my_inbox.get_message(i)
    # Return (from_number, time_arrived, text_of_sms) for message[i]
    # Also change its state to "has been viewed".
    # If there is no message at position i, return None

my_inbox.delete(i)      # Delete the message at index i
my_inbox.clear()        # Delete all messages from inbox
```

Write the class, create a message store object, write tests for these methods, and implement the methods.

## 11.2 Classes and Objects — Digging a little deeper

### 11.2.1 Rectangles

Let's say that we want a class to represent a rectangle which is located somewhere in the XY plane. The question is, what information do we have to provide in order to specify such a rectangle? To keep things simple, assume that the rectangle is oriented either vertically or horizontally, never at an angle.

There are a few possibilities: we could specify the center of the rectangle (two coordinates) and its size (width and height); or we could specify one of the corners and the size; or we could specify two opposing corners. A conventional choice is to specify the upper-left corner of the rectangle, and the size.

Again, we'll define a new class, and provide it with an initializer and a string converter method:

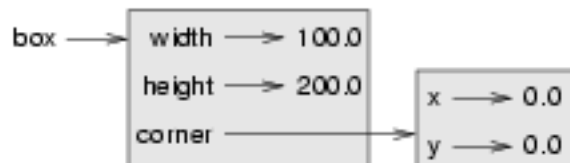
```
1 class Rectangle:
2     """ A class to manufacture rectangle objects """
3
4     def __init__(self, posn, w, h):
5         """ Initialize rectangle at posn, with width w, height h """
6         self.corner = posn
7         self.width = w
8         self.height = h
9
10    def __str__(self):
11        return "{0}, {1}, {2}"
12            .format(self.corner, self.width, self.height)
13
14    box = Rectangle(Point(0, 0), 100, 200)
15    bomb = Rectangle(Point(100, 80), 5, 10)    # In my video game
16    print("box: ", box)
17    print("bomb: ", bomb)
```

To specify the upper-left corner, we have embedded a `Point` object (as we used it in the previous chapter) within our new `Rectangle` object! We create two new `Rectangle` objects, and then print them producing:

```
box: ((0, 0), 100, 200)
bomb: ((100, 80), 5, 10)
```

The dot operator composes. The expression `box.corner.x` means, “Go to the object that `box` refers to and select its attribute named `corner`, then go to that object and select its attribute named `x`”.

The figure shows the state of this object:



### 11.2.2 Objects are mutable

We can change the state of an object by making an assignment to one of its attributes. For example, to grow the size of a rectangle without changing its position, we could modify the values of `width` and `height`:

```
box.width += 50
box.height += 100
```

Of course, we'd probably like to provide a method to encapsulate this inside the class. We will also provide another method to move the position of the rectangle elsewhere:

```
1 class Rectangle:
2     # ...
3
4     def grow(self, delta_width, delta_height):
5         """ Grow (or shrink) this object by the deltas """
6         self.width += delta_width
7         self.height += delta_height
8
9     def move(self, dx, dy):
10        """ Move this object by the deltas """
11        self.corner.x += dx
12        self.corner.y += dy
```

Let us try this:

```
>>> r = Rectangle(Point(10,5), 100, 50)
>>> print(r)
((10, 5), 100, 50)
>>> r.grow(25, -10)
>>> print(r)
((10, 5), 125, 40)
>>> r.move(-10, 10)
print(r)
((0, 15), 125, 40)
```

### 11.2.3 Sameness

The meaning of the word “same” seems perfectly clear until we give it some thought, and then we realize there is more to it than we initially expected.

For example, if we say, “Alice and Bob have the same car”, we mean that her car and his are the same make and model, but that they are two different cars. If we say, “Alice and Bob have the same mother”, we mean that her mother and his are the same person.

When we talk about objects, there is a similar ambiguity. For example, if two `Point`s are the same, does that mean they contain the same data (coordinates) or that they are actually the same object?

We’ve already seen the `is` operator in the chapter on lists, where we talked about aliases: it allows us to find out if two references refer to the same object:

```
>>> p1 = Point(3, 4)
>>> p2 = Point(3, 4)
>>> p1 is p2
False
```

Even though `p1` and `p2` contain the same coordinates, they are not the same object. If we assign `p1` to `p3`, then the two variables are aliases of the same object:

```
>>> p3 = p1
>>> p1 is p3
True
```

This type of equality is called **shallow equality** because it compares only the references, not the contents of the objects.

To compare the contents of the objects — **deep equality** — we can write a function called `same_coordinates`:

```
1 def same_coordinates(p1, p2):
2     return (p1.x == p2.x) and (p1.y == p2.y)
```

Now if we create two different objects that contain the same data, we can use `same_point` to find out if they represent points with the same coordinates.

```
>>> p1 = Point(3, 4)
>>> p2 = Point(3, 4)
>>> same_coordinates(p1, p2)
True
```

Of course, if the two variables refer to the same object, they have both shallow and deep equality.

---

### Beware of ==

“When I use a word,” Humpty Dumpty said, in a rather scornful tone, “it means just what I choose it to mean — neither more nor less.” *Alice in Wonderland*

Python has a powerful feature that allows a designer of a class to decide what an operation like `==` or `<` should mean. (We’ve just shown how we can control how our own objects are converted to strings, so we’ve already made a start!) We’ll cover more detail later. But sometimes the implementors will attach shallow equality semantics, and sometimes deep equality, as shown in this little experiment:

```
1 p = Point(4, 2)
2 s = Point(4, 2)
3 print("== on Points returns", p == s)
4 # By default, == on Point objects does a shallow equality test
5
6 a = [2,3]
7 b = [2,3]
8 print("== on lists returns", a == b)
9 # But by default, == does a deep equality test on lists
```

This outputs:

```
== on Points returns False
== on lists returns True
```

So we conclude that even though the two lists (or tuples, etc.) are distinct objects with different memory addresses, for lists the `==` operator tests for deep equality, while in the case of points it makes a shallow test.

---

### 11.2.4 Copying

Aliasing can make a program difficult to read because changes made in one place might have unexpected effects in another place. It is hard to keep track of all the variables that might refer to a given object.

Copying an object is often an alternative to aliasing. The `copy` module contains a function called `copy` that can duplicate any object:

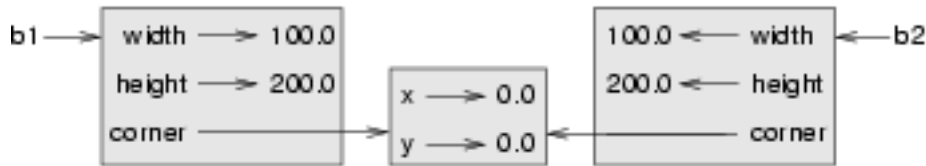
```
>>> import copy
>>> p1 = Point(3, 4)
>>> p2 = copy.copy(p1)
>>> p1 is p2
False
>>> same_coordinates(p1, p2)
True
```

Once we import the `copy` module, we can use the `copy` function to make a new `Point`. `p1` and `p2` are not the same point, but they contain the same data.

To copy a simple object like a `Point`, which doesn't contain any embedded objects, `copy` is sufficient. This is called **shallow copying**.

For something like a `Rectangle`, which contains a reference to a `Point`, `copy` doesn't do quite the right thing. It copies the reference to the `Point` object, so both the old `Rectangle` and the new one refer to a single `Point`.

If we create a box, `b1`, in the usual way and then make a copy, `b2`, using `copy`, the resulting state diagram looks like this:



This is almost certainly not what we want. In this case, invoking `grow` on one of the `Rectangle` objects would not affect the other, but invoking `move` on either would affect both! This behavior is confusing and error-prone. The shallow copy has created an alias to the `Point` that represents the corner.

Fortunately, the `copy` module contains a function named `deepcopy` that copies not only the object but also any embedded objects. It won't be surprising to learn that this operation is called a **deep copy**.

```
>>> b2 = copy.deepcopy(b1)
```

Now `b1` and `b2` are completely separate objects.

## 11.2.5 Glossary

**deep copy** To copy the contents of an object as well as any embedded objects, and any objects embedded in them, and so on; implemented by the `deepcopy` function in the `copy` module.

**deep equality** Equality of values, or two references that point to objects that have the same value.

**shallow copy** To copy the contents of an object, including any references to embedded objects; implemented by the `copy` function in the `copy` module.

**shallow equality** Equality of references, or two references that point to the same object.

## 11.2.6 Exercises

1. Add a method `area` to the `Rectangle` class that returns the area of any instance:

```
r = Rectangle(Point(0, 0), 10, 5)
test(r.area() == 50)
```

2. Write a `perimeter` method in the `Rectangle` class so that we can find the perimeter of any rectangle instance:

```
r = Rectangle(Point(0, 0), 10, 5)
test(r.perimeter() == 30)
```

3. Write a `flip` method in the `Rectangle` class that swaps the width and the height of any rectangle instance:

```
r = Rectangle(Point(100, 50), 10, 5)
test(r.width == 10 and r.height == 5)
r.flip()
test(r.width == 5 and r.height == 10)
```

4. Write a new method in the `Rectangle` class to test if a `Point` falls within the rectangle. For this exercise, assume that a rectangle at (0,0) with width 10 and height 5 has *open* upper bounds on the width and height, i.e. it stretches in the x direction from [0 to 10), where 0 is included but 10 is excluded, and from [0 to 5) in the y direction. So it does not contain the point (10, 2). These tests should pass:

```
r = Rectangle(Point(0, 0), 10, 5)
test(r.contains(Point(0, 0)))
test(r.contains(Point(3, 3)))
test(not r.contains(Point(3, 7)))
test(not r.contains(Point(3, 5)))
test(r.contains(Point(3, 4.99999)))
test(not r.contains(Point(-3, -3)))
```

5. In games, we often put a rectangular “bounding box” around our sprites. (A sprite is an object that can move about in the game, as we will see shortly.) We can then do *collision detection* between, say, bombs and spaceships, by comparing whether their rectangles overlap anywhere.

Write a function to determine whether two rectangles collide. *Hint: this might be quite a tough exercise! Think carefully about all the cases before you code.*

## 11.3 Even more OOP

### 11.3.1 MyTime

As another example of a user-defined type, we’ll define a class called `MyTime` that records the time of day. We’ll provide an `__init__` method to ensure that every instance is created with appropriate attributes and initialization. The class definition looks like this:

```
1 class MyTime:
2
3     def __init__(self, hrs=0, mins=0, secs=0):
4         """ Create a MyTime object initialized to hrs, mins, secs """
5         self.hours = hrs
6         self.minutes = mins
7         self.seconds = secs
```

We can instantiate a new `MyTime` object:

```
1 tim1 = MyTime(11, 59, 30)
```

The state diagram for the object looks like this:



We'll leave it as an exercise for the readers to add a `__str__` method so that `MyTime` objects can print themselves decently.

### 11.3.2 Pure functions

In the next few sections, we'll write two versions of a function called `add_time`, which calculates the sum of two `MyTime` objects. They will demonstrate two kinds of functions: pure functions and modifiers.

The following is a rough version of `add_time`:

```
1 def add_time(t1, t2):
2     h = t1.hours + t2.hours
3     m = t1.minutes + t2.minutes
4     s = t1.seconds + t2.seconds
5     sum_t = MyTime(h, m, s)
6     return sum_t
```

The function creates a new `MyTime` object and returns a reference to the new object. This is called a **pure function** because it does not modify any of the objects passed to it as parameters and it has no side effects, such as updating global variables, displaying a value, or getting user input.

Here is an example of how to use this function. We'll create two `MyTime` objects: `current_time`, which contains the current time; and `bread_time`, which contains the amount of time it takes for a breadmaker to make bread. Then we'll use `add_time` to figure out when the bread will be done.

```
>>> current_time = MyTime(9, 14, 30)
>>> bread_time = MyTime(3, 35, 0)
>>> done_time = add_time(current_time, bread_time)
>>> print(done_time)
12:49:30
```

The output of this program is 12:49:30, which is correct. On the other hand, there are cases where the result is not correct. Can you think of one?

The problem is that this function does not deal with cases where the number of seconds or minutes adds up to more than sixty. When that happens, we have to carry the extra seconds into the minutes column or the extra minutes into the hours column.

Here's a better version of the function:

```
1 def add_time(t1, t2):
2
3     h = t1.hours + t2.hours
4     m = t1.minutes + t2.minutes
5     s = t1.seconds + t2.seconds
6
7     if s >= 60:
8         s -= 60
```

(continues on next page)



(continued from previous page)

```
9         m += 1
10
11     if m >= 60:
12         m -= 60
13         h += 1
14
15     sum_t = MyTime(h, m, s)
16     return sum_t
```

This function is starting to get bigger, and still doesn't work for all possible cases. Later we will suggest an alternative approach that yields better code.

### 11.3.3 Modifiers

There are times when it is useful for a function to modify one or more of the objects it gets as parameters. Usually, the caller keeps a reference to the objects it passes, so any changes the function makes are visible to the caller. Functions that work this way are called **modifiers**.

`increment`, which adds a given number of seconds to a `MyTime` object, would be written most naturally as a modifier. A rough draft of the function looks like this:

```
1 def increment(t, secs):
2     t.seconds += secs
3
4     if t.seconds >= 60:
5         t.seconds -= 60
6         t.minutes += 1
7
8     if t.minutes >= 60:
9         t.minutes -= 60
10        t.hours += 1
```

The first line performs the basic operation; the remainder deals with the special cases we saw before.

Is this function correct? What happens if the parameter `seconds` is much greater than sixty? In that case, it is not enough to carry once; we have to keep doing it until `seconds` is less than sixty. One solution is to replace the `if` statements with `while` statements:

```
1 def increment(t, seconds):
2     t.seconds += seconds
3
4     while t.seconds >= 60:
5         t.seconds -= 60
6         t.minutes += 1
7
8     while t.minutes >= 60:
9         t.minutes -= 60
10        t.hours += 1
```

This function is now correct when `seconds` is not negative, and when `hours` does not exceed 23, but it is not a particularly good solution.

### 11.3.4 Converting `increment` to a method

Once again, OOP programmers would prefer to put functions that work with `MyTime` objects into the `MyTime` class, so let's convert `increment` to a method. To save space, we will leave out previously defined methods, but you should keep them in your version:

```

1 class MyTime:
2     # Previous method definitions here...
3
4     def increment(self, seconds):
5         self.seconds += seconds
6
7         while self.seconds >= 60:
8             self.seconds -= 60
9             self.minutes += 1
10
11        while self.minutes >= 60:
12            self.minutes -= 60
13            self.hours += 1

```

The transformation is purely mechanical — we move the definition into the class definition and (optionally) change the name of the first parameter to `self`, to fit with Python style conventions.

Now we can invoke `increment` using the syntax for invoking a method.

```

1 current_time.increment(500)

```

Again, the object on which the method is invoked gets assigned to the first parameter, `self`. The second parameter, `seconds` gets the value 500.

### 11.3.5 An “Aha!” insight

Often a high-level insight into the problem can make the programming much easier.

In this case, the insight is that a `MyTime` object is really a three-digit number in base 60! The second component is the ones column, the minute component is the sixties column, and the hour component is the thirty-six hundreds column.

When we wrote `add_time` and `increment`, we were effectively doing addition in base 60, which is why we had to carry from one column to the next.

This observation suggests another approach to the whole problem — we can convert a `MyTime` object into a single number and take advantage of the fact that the computer knows how to do arithmetic with numbers. The following method is added to the `MyTime` class to convert any instance into a corresponding number of seconds:

```

1 class MyTime:
2     # ...
3
4     def to_seconds(self):
5         """ Return the number of seconds represented
6             by this instance
7         """
8         return self.hours * 3600 + self.minutes * 60 + self.seconds

```

Now, all we need is a way to convert from an integer back to a `MyTime` object. Supposing we have `tsecs` seconds, some integer division and mod operators can do this for us:

```
1 hrs = tsecs // 3600
2 leftoversecs = tsecs % 3600
3 mins = leftoversecs // 60
4 secs = leftoversecs % 60
```

You might have to think a bit to convince yourself that this technique to convert from one base to another is correct.

In OOP we're really trying to wrap together the data and the operations that apply to it. So we'd like to have this logic inside the `MyTime` class. A good solution is to rewrite the class initializer so that it can cope with initial values of seconds or minutes that are outside the **normalized** values. (A normalized time would be something like 3 hours 12 minutes and 20 seconds. The same time, but unnormalized could be 2 hours 70 minutes and 140 seconds.)

Let's rewrite a more powerful initializer for `MyTime`:

```
1 class MyTime:
2     # ...
3
4     def __init__(self, hrs=0, mins=0, secs=0):
5         """ Create a new MyTime object initialized to hrs, mins, secs.
6             The values of mins and secs may be outside the range 0-59,
7             but the resulting MyTime object will be normalized.
8         """
9
10        # Calculate total seconds to represent
11        totalsecs = hrs*3600 + mins*60 + secs
12        self.hours = totalsecs // 3600          # Split in h, m, s
13        leftoversecs = totalsecs % 3600
14        self.minutes = leftoversecs // 60
15        self.seconds = leftoversecs % 60
```

Now we can rewrite `add_time` like this:

```
1 def add_time(t1, t2):
2     secs = t1.to_seconds() + t2.to_seconds()
3     return MyTime(0, 0, secs)
```

This version is much shorter than the original, and it is much easier to demonstrate or reason that it is correct.

### 11.3.6 Generalization

In some ways, converting from base 60 to base 10 and back is harder than just dealing with times. Base conversion is more abstract; our intuition for dealing with times is better.

But if we have the insight to treat times as base 60 numbers and make the investment of writing the conversions, we get a program that is shorter, easier to read and debug, and more reliable.

It is also easier to add features later. For example, imagine subtracting two `MyTime` objects to find the duration between them. The naive approach would be to implement subtraction with borrowing. Using the conversion functions would be easier and more likely to be correct.

Ironically, sometimes making a problem harder (or more general) makes the programming easier, because there are fewer special cases and fewer opportunities for error.

---

#### Specialization versus Generalization

Computer Scientists are generally fond of specializing their types, while mathematicians often take the opposite approach, and generalize everything.

What do we mean by this?

If we ask a mathematician to solve a problem involving weekdays, days of the century, playing cards, time, or dominoes, their most likely response is to observe that all these objects can be represented by integers. Playing cards, for example, can be numbered from 0 to 51. Days within the century can be numbered. Mathematicians will say “*These things are enumerable — the elements can be uniquely numbered (and we can reverse this numbering to get back to the original concept). So let’s number them, and confine our thinking to integers. Luckily, we have powerful techniques and a good understanding of integers, and so our abstractions — the way we tackle and simplify these problems — is to try to reduce them to problems about integers.*”

Computer Scientists tend to do the opposite. We will argue that there are many integer operations that are simply not meaningful for dominoes, or for days of the century. So we’ll often define new specialized types, like `MyTime`, because we can restrict, control, and specialize the operations that are possible. Object-oriented programming is particularly popular because it gives us a good way to bundle methods and specialized data into a new type.

Both approaches are powerful problem-solving techniques. Often it may help to try to think about the problem from both points of view — “*What would happen if I tried to reduce everything to very few primitive types?*”, versus “*What would happen if this thing had its own specialized type?*”

---

### 11.3.7 Another example

The `after` function should compare two times, and tell us whether the first time is strictly after the second, e.g.

```
>>> t1 = MyTime(10, 55, 12)
>>> t2 = MyTime(10, 48, 22)
>>> after(t1, t2)           # Is t1 after t2?
True
```

This is slightly more complicated because it operates on two `MyTime` objects, not just one. But we’d prefer to write it as a method anyway — in this case, a method on the first argument:

```
1 class MyTime:
2     # Previous method definitions here...
3
4     def after(self, time2):
5         """ Return True if I am strictly greater than time2 """
6         if self.hours > time2.hours:
7             return True
8         if self.hours < time2.hours:
9             return False
10
11        if self.minutes > time2.minutes:
12            return True
13        if self.minutes < time2.minutes:
14            return False
15        if self.seconds > time2.seconds:
16            return True
17
18        return False
```

We invoke this method on one object and pass the other as an argument:

```
1 if current_time.after(done_time):
2     print("The bread will be done before it starts!")
```

We can almost read the invocation like English: If the current time is after the done time, then...

The logic of the `if` statements deserve special attention here. Lines 11-18 will only be reached if the two hour fields are the same. Similarly, the test at line 16 is only executed if both times have the same hours and the same minutes.

Could we make this easier by using our “Aha!” insight and extra work from earlier, and reducing both times to integers? Yes, with spectacular results!

```
1 class MyTime:
2     # Previous method definitions here...
3
4     def after(self, time2):
5         """ Return True if I am strictly greater than time2 """
6         return self.to_seconds() > time2.to_seconds()
```

This is a great way to code this: if we want to tell if the first time is after the second time, turn them both into integers and compare the integers.

### 11.3.8 Operator overloading

Some languages, including Python, make it possible to have different meanings for the same operator when applied to different types. For example, `+` in Python means quite different things for integers and for strings. This feature is called **operator overloading**.

It is especially useful when programmers can also overload the operators for their own user-defined types.

For example, to override the addition operator `+`, we can provide a method named `__add__`:

```
1 class MyTime:
2     # Previously defined methods here...
3
4     def __add__(self, other):
5         return MyTime(0, 0, self.to_seconds() + other.to_seconds())
```

As usual, the first parameter is the object on which the method is invoked. The second parameter is conveniently named `other` to distinguish it from `self`. To add two `MyTime` objects, we create and return a new `MyTime` object that contains their sum.

Now, when we apply the `+` operator to `MyTime` objects, Python invokes the `__add__` method that we have written:

```
>>> t1 = MyTime(1, 15, 42)
>>> t2 = MyTime(3, 50, 30)
>>> t3 = t1 + t2
>>> print(t3)
05:06:12
```

The expression `t1 + t2` is equivalent to `t1.__add__(t2)`, but obviously more elegant. As an exercise, add a method `__sub__(self, other)` that overloads the subtraction operator, and try it out.

For the next couple of exercises we'll go back to the `Point` class defined in our first chapter about objects, and overload some of its operators. Firstly, adding two points adds their respective (x, y) coordinates:

```
1 class Point:
2     # Previously defined methods here...
3
4     def __add__(self, other):
5         return Point(self.x + other.x, self.y + other.y)
```

There are several ways to override the behavior of the multiplication operator: by defining a method named `__mul__`, or `__rmul__`, or both.

If the left operand of `*` is a `Point`, Python invokes `__mul__`, which assumes that the other operand is also a `Point`. It computes the **dot product** of the two `Points`, defined according to the rules of linear algebra:

```
1 def __mul__(self, other):
2     return self.x * other.x + self.y * other.y
```

If the left operand of `*` is a primitive type and the right operand is a `Point`, Python invokes `__rmul__`, which performs **scalar multiplication**:

```
1 def __rmul__(self, other):
2     return Point(other * self.x, other * self.y)
```

The result is a new `Point` whose coordinates are a multiple of the original coordinates. If `other` is a type that cannot be multiplied by a floating-point number, then `__rmul__` will yield an error.

This example demonstrates both kinds of multiplication:

```
>>> p1 = Point(3, 4)
>>> p2 = Point(5, 7)
>>> print(p1 * p2)
43
>>> print(2 * p2)
(10, 14)
```

What happens if we try to evaluate `p2 * 2`? Since the first parameter is a `Point`, Python invokes `__mul__` with 2 as the second argument. Inside `__mul__`, the program tries to access the `x` coordinate of `other`, which fails because an integer has no attributes:

```
>>> print(p2 * 2)
AttributeError: 'int' object has no attribute 'x'
```

Unfortunately, the error message is a bit opaque. This example demonstrates some of the difficulties of object-oriented programming. Sometimes it is hard enough just to figure out what code is running.

### 11.3.9 Polymorphism

Most of the methods we have written only work for a specific type. When we create a new object, we write methods that operate on that type.

But there are certain operations that we will want to apply to many types, such as the arithmetic operations in the previous sections. If many types support the same set of operations, we can write functions that work on any of those types.

For example, the `multadd` operation (which is common in linear algebra) takes three parameters; it multiplies the first two and then adds the third. We can write it in Python like this:

```
1 def multadd(x, y, z):
2     return x * y + z
```

This function will work for any values of `x` and `y` that can be multiplied and for any value of `z` that can be added to the product.

We can invoke it with numeric values:

```
>>> multadd(3, 2, 1)
7
```

Or with `Points`:

```
>>> p1 = Point(3, 4)
>>> p2 = Point(5, 7)
>>> print(multadd (2, p1, p2))
(11, 15)
>>> print(multadd (p1, p2, 1))
44
```

In the first case, the `Point` is multiplied by a scalar and then added to another `Point`. In the second case, the dot product yields a numeric value, so the third parameter also has to be a numeric value.

A function like this that can take arguments with different types is called **polymorphic**.

As another example, consider the function `front_and_back`, which prints a list twice, forward and backward:

```
1 def front_and_back(front):
2     import copy
3     back = copy.copy(front)
4     back.reverse()
5     print(str(front) + str(back))
```

Because the `reverse` method is a modifier, we make a copy of the list before reversing it. That way, this function doesn't modify the list it gets as a parameter.

Here's an example that applies `front_and_back` to a list:

```
>>> my_list = [1, 2, 3, 4]
>>> front_and_back(my_list)
[1, 2, 3, 4][4, 3, 2, 1]
```

Of course, we intended to apply this function to lists, so it is not surprising that it works. What would be surprising is if we could apply it to a `Point`.

To determine whether a function can be applied to a new type, we apply Python's fundamental rule of polymorphism, called the **duck typing rule**: *If all of the operations inside the function can be applied to the type, the function can be applied to the type.* The operations in the `front_and_back` function include `copy`, `reverse`, and `print`.

Not all programming languages define polymorphism in this way. Look up *duck typing*, and see if you can figure out why it has this name.

`copy` works on any object, and we have already written a `__str__` method for `Point` objects, so all we need is a `reverse` method in the `Point` class:

```
1 def reverse(self):
2     (self.x, self.y) = (self.y, self.x)
```

Then we can pass `Points` to `front_and_back`:

```
>>> p = Point(3, 4)
>>> front_and_back(p)
(3, 4) (4, 3)
```

The most interesting polymorphism is the unintentional kind, where we discover that a function we have already written can be applied to a type for which we never planned.

### 11.3.10 Glossary

**dot product** An operation defined in linear algebra that multiplies two `Points` and yields a numeric value.

**functional programming style** A style of program design in which the majority of functions are pure.

**modifier** A function or method that changes one or more of the objects it receives as parameters. Most modifier functions are void (do not return a value).

**normalized** Data is said to be normalized if it fits into some reduced range or set of rules. We usually normalize our angles to values in the range [0..360). We normalize minutes and seconds to be values in the range [0..60). And we'd be surprised if the local store advertised its cold drinks at "One dollar, two hundred and fifty cents".

**operator overloading** Extending built-in operators ( +, −, \*, >, <, etc.) so that they do different things for different types of arguments. We've seen early in the book how + is overloaded for numbers and strings, and here we've shown how to further overload it for user-defined types.

**polymorphic** A function that can operate on more than one type. Notice the subtle distinction: overloading has different functions (all with the same name) for different types, whereas a polymorphic function is a single function that can work for a range of types.

**pure function** A function that does not modify any of the objects it receives as parameters. Most pure functions are fruitful rather than void.

**scalar multiplication** An operation defined in linear algebra that multiplies each of the coordinates of a `Point` by a numeric value.

## 11.3.11 Exercises

1. Write a Boolean function `between` that takes two `MyTime` objects, `t1` and `t2`, as arguments, and returns `True` if the invoking object falls between the two times. Assume `t1 <= t2`, and make the test closed at the lower bound and open at the upper bound, i.e. return `True` if `t1 <= obj < t2`.
2. Turn the above function into a method in the `MyTime` class.
3. Overload the necessary operator(s) so that instead of having to write

```
if t1.after(t2): ...
```

we can use the more convenient

```
if t1 > t2: ...
```

4. Rewrite `increment` as a method that uses our "Aha" insight.
5. Create some test cases for the `increment` method. Consider specifically the case where the number of seconds to add to the time is negative. Fix up `increment` so that it handles this case if it does not do so already. (You may assume that you will never subtract more seconds than are in the time object.)
6. Can physical time be negative, or must time always move in the forward direction? Some serious physicists think this is not such a dumb question. See what you can find on the Internet about this.

## 11.4 Collections of objects

### 11.4.1 Composition

By now, we have seen several examples of composition. One of the first examples was using a method invocation as part of an expression. Another example is the nested structure of statements; we can put an `if` statement within a



`while` loop, within another `if` statement, and so on.

Having seen this pattern, and having learned about lists and objects, we should not be surprised to learn that we can create lists of objects. We can also create objects that contain lists (as attributes); we can create lists that contain lists; we can create objects that contain objects; and so on.

In this chapter and the next, we will look at some examples of these combinations, using `Card` objects as an example.

### 11.4.2 Card objects

If you are not familiar with common playing cards, now would be a good time to get a deck, or else this chapter might not make much sense. There are fifty-two cards in a deck, each of which belongs to one of four suits and one of thirteen ranks. The suits are Spades, Hearts, Diamonds, and Clubs (in descending order in bridge). The ranks are Ace, 2, 3, 4, 5, 6, 7, 8, 9, 10, Jack, Queen, and King. Depending on the game that we are playing, the rank of Ace may be higher than King or lower than 2. The rank is sometimes called the face-value of the card.

If we want to define a new object to represent a playing card, it is obvious what the attributes should be: `rank` and `suit`. It is not as obvious what type the attributes should be. One possibility is to use strings containing words like "Spade" for suits and "Queen" for ranks. One problem with this implementation is that it would not be easy to compare cards to see which had a higher rank or suit.

An alternative is to use integers to **encode** the ranks and suits. By encode, we do not mean what some people think, which is to encrypt or translate into a secret code. What a computer scientist means by encode is to define a mapping between a sequence of numbers and the items I want to represent. For example:

```
Spades    -->  3
Hearts    -->  2
Diamonds  -->  1
Clubs     -->  0
```

An obvious feature of this mapping is that the suits map to integers in order, so we can compare suits by comparing integers. The mapping for ranks is fairly obvious; each of the numerical ranks maps to the corresponding integer, and for face cards:

```
Jack      --> 11
Queen     --> 12
King      --> 13
```

The reason we are using mathematical notation for these mappings is that they are not part of the Python program. They are part of the program design, but they never appear explicitly in the code. The class definition for the `Card` type looks like this:

```
1 class Card:
2     def __init__(self, suit=0, rank=0):
3         self.suit = suit
4         self.rank = rank
```

As usual, we provide an initialization method that takes an optional parameter for each attribute.

To create some objects, representing say the 3 of Clubs and the Jack of Diamonds, use these commands:

```
1 three_of_clubs = Card(0, 3)
2 card1 = Card(1, 11)
```

In the first case above, for example, the first argument, 0, represents the suit Clubs.

---

Save this code for later use ...

In the next chapter we assume that we have save the `Cards` class, and the upcoming `Deck` class in a file called `Cards.py`.

---

### 11.4.3 Class attributes and the `__str__` method

In order to print `Card` objects in a way that people can easily read, we want to map the integer codes onto words. A natural way to do that is with lists of strings. We assign these lists to **class attributes** at the top of the class definition:

```
1 class Card:
2     suits = ["Clubs", "Diamonds", "Hearts", "Spades"]
3     ranks = ["narf", "Ace", "2", "3", "4", "5", "6", "7",
4             "8", "9", "10", "Jack", "Queen", "King"]
5
6     def __init__(self, suit=0, rank=0):
7         self.suit = suit
8         self.rank = rank
9
10    def __str__(self):
11        return (self.ranks[self.rank] + " of " + self.suits[self.suit])
```

A class attribute is defined outside of any method, and it can be accessed from any of the methods in the class.

Inside `__str__`, we can use `suits` and `ranks` to map the numerical values of `suit` and `rank` to strings. For example, the expression `self.suits[self.suit]` means use the attribute `suit` from the object `self` as an index into the class attribute named `suits`, and select the appropriate string.

The reason for the "narf" in the first element in `ranks` is to act as a place keeper for the zero-eth element of the list, which will never be used. The only valid ranks are 1 to 13. This wasted item is not entirely necessary. We could have started at 0, as usual, but it is less confusing to encode the rank 2 as integer 2, 3 as 3, and so on.

With the methods we have so far, we can create and print cards:

```
>>> card1 = Card(1, 11)
>>> print(card1)
Jack of Diamonds
```

Class attributes like `suits` are shared by all `Card` objects. The advantage of this is that we can use any `Card` object to access the class attributes:

```
>>> card2 = Card(1, 3)
>>> print(card2)
3 of Diamonds
>>> print(card2.suits[1])
Diamonds
```

Because every `Card` instance references the same class attribute, we have an aliasing situation. The disadvantage is that if we modify a class attribute, it affects every instance of the class. For example, if we decide that Jack of Diamonds should really be called Jack of Swirly Whales, we could do this:

```
>>> card1.suits[1] = "Swirly Whales"
>>> print(card1)
Jack of Swirly Whales
```

The problem is that *all* of the Diamonds just became Swirly Whales:

```
>>> print(card2)
3 of Swirly Whales
```

It is usually not a good idea to modify class attributes.

### 11.4.4 Comparing cards

For primitive types, there are six relational operators (<, >, ==, etc.) that compare values and determine when one is greater than, less than, or equal to another. If we want our own types to be comparable using the syntax of these relational operators, we need to define six corresponding special methods in our class.

We'd like to start with a single method named `cmp` that houses the logic of ordering. By convention, a comparison method takes two parameters, `self` and `other`, and returns 1 if the first object is greater, -1 if the second object is greater, and 0 if they are equal to each other.

Some types are completely ordered, which means that we can compare any two elements and tell which is bigger. For example, the integers and the floating-point numbers are completely ordered. Some types are unordered, which means that there is no meaningful way to say that one element is bigger than another. For example, the fruits are unordered, which is why we cannot compare apples and oranges, and we cannot meaningfully order a collection of images, or a collection of cellphones.

Playing cards are partially ordered, which means that sometimes we can compare cards and sometimes not. For example, we know that the 3 of Clubs is higher than the 2 of Clubs, and the 3 of Diamonds is higher than the 3 of Clubs. But which is better, the 3 of Clubs or the 2 of Diamonds? One has a higher rank, but the other has a higher suit.

In order to make cards comparable, we have to decide which is more important, rank or suit. To be honest, the choice is arbitrary. For the sake of choosing, we will say that suit is more important, because a new deck of cards comes sorted with all the Clubs together, followed by all the Diamonds, and so on.

With that decided, we can write `cmp`:

```
1 def cmp(self, other):
2     # Check the suits
3     if self.suit > other.suit: return 1
4     if self.suit < other.suit: return -1
5     # Suits are the same... check ranks
6     if self.rank > other.rank: return 1
7     if self.rank < other.rank: return -1
8     # Ranks are the same... it's a tie
9     return 0
```

In this ordering, Aces appear lower than Deuces (2s).

Now, we can define the six special methods that do the overloading of each of the relational operators for us:

```
1 def __eq__(self, other):
2     return self.cmp(other) == 0
3
4 def __le__(self, other):
5     return self.cmp(other) <= 0
6
7 def __ge__(self, other):
8     return self.cmp(other) >= 0
9
10 def __gt__(self, other):
11     return self.cmp(other) > 0
12
```

(continues on next page)

(continued from previous page)

```
13 def __lt__(self, other):
14     return self.cmp(other) < 0
15
16 def __ne__(self, other):
17     return self.cmp(other) != 0
```

With this machinery in place, the relational operators now work as we'd like them to:

```
>>> card1 = Card(1, 11)
>>> card2 = Card(1, 3)
>>> card3 = Card(1, 11)
>>> card1 < card2
False
>>> card1 == card3
True
```

### 11.4.5 Decks

Now that we have objects to represent Cards, the next logical step is to define a class to represent a Deck. Of course, a deck is made up of cards, so each Deck object will contain a list of cards as an attribute. Many card games will need at least two different decks — a red deck and a blue deck.

The following is a class definition for the Deck class. The initialization method creates the attribute `cards` and generates the standard pack of fifty-two cards:

```
1 class Deck:
2     def __init__(self):
3         self.cards = []
4         for suit in range(4):
5             for rank in range(1, 14):
6                 self.cards.append(Card(suit, rank))
```

The easiest way to populate the deck is with a nested loop. The outer loop enumerates the suits from 0 to 3. The inner loop enumerates the ranks from 1 to 13. Since the outer loop iterates four times, and the inner loop iterates thirteen times, the total number of times the body is executed is fifty-two (thirteen times four). Each iteration creates a new instance of Card with the current suit and rank, and appends that card to the `cards` list.

With this in place, we can instantiate some decks:

```
1 red_deck = Deck()
2 blue_deck = Deck()
```

### 11.4.6 Printing the deck

As usual, when we define a new type we want a method that prints the contents of an instance. To print a Deck, we traverse the list and print each Card:

```
1 class Deck:
2     ...
3     def print_deck(self):
4         for card in self.cards:
5             print(card)
```

Here, and from now on, the ellipsis ( . . . ) indicates that we have omitted the other methods in the class.

As an alternative to `print_deck`, we could write a `__str__` method for the `Deck` class. The advantage of `__str__` is that it is more flexible. Rather than just printing the contents of the object, it generates a string representation that other parts of the program can manipulate before printing, or store for later use.

Here is a version of `__str__` that returns a string representation of a `Deck`. To add a bit of pizzazz, it arranges the cards in a cascade where each card is indented one space more than the previous card:

```
1 class Deck:
2     ...
3     def __str__(self):
4         s = ""
5         for i in range(len(self.cards)):
6             s = s + " " * i + str(self.cards[i]) + "\n"
7         return s
```

This example demonstrates several features. First, instead of traversing `self.cards` and assigning each card to a variable, we are using `i` as a loop variable and an index into the list of cards.

Second, we are using the string multiplication operator to indent each card by one more space than the last. The expression `" " * i` yields a number of spaces equal to the current value of `i`.

Third, instead of using the `print` command to print the cards, we use the `str` function. Passing an object as an argument to `str` is equivalent to invoking the `__str__` method on the object.

Finally, we are using the variable `s` as an **accumulator**. Initially, `s` is the empty string. Each time through the loop, a new string is generated and concatenated with the old value of `s` to get the new value. When the loop ends, `s` contains the complete string representation of the `Deck`, which looks like this:

```
>>> red_deck = Deck()
>>> print(red_deck)
Ace of Clubs
 2 of Clubs
 3 of Clubs
 4 of Clubs
 5 of Clubs
 6 of Clubs
 7 of Clubs
 8 of Clubs
 9 of Clubs
10 of Clubs
 Jack of Clubs
  Queen of Clubs
  King of Clubs
   Ace of Diamonds
    2 of Diamonds
    ...
```

And so on. Even though the result appears on 52 lines, it is one long string that contains newlines.

### 11.4.7 Shuffling the deck

If a deck is perfectly shuffled, then any card is equally likely to appear anywhere in the deck, and any location in the deck is equally likely to contain any card.

To shuffle the deck, we will use the `randrange` function from the `random` module. With two integer arguments, `a` and `b`, `randrange` chooses a random integer in the range `a <= x < b`. Since the upper bound is strictly less than `b`, we can use the length of a list as the second parameter, and we are guaranteed to get a legal index. For example, if

`rng` has already been instantiated as a random number source, this expression chooses the index of a random card in a deck:

```
1 rng.randrange(0, len(self.cards))
```

An easy way to shuffle the deck is by traversing the cards and swapping each card with a randomly chosen one. It is possible that the card will be swapped with itself, but that is fine. In fact, if we precluded that possibility, the order of the cards would be less than entirely random:

```
1 class Deck:
2     ...
3     def shuffle(self):
4         import random
5         rng = random.Random()           # Create a random generator
6         num_cards = len(self.cards)
7         for i in range(num_cards):
8             j = rng.randrange(i, num_cards)
9             (self.cards[i], self.cards[j]) = (self.cards[j], self.cards[i])
```

Rather than assume that there are fifty-two cards in the deck, we get the actual length of the list and store it in `num_cards`.

For each card in the deck, we choose a random card from among the cards that haven't been shuffled yet. Then we swap the current card (`i`) with the selected card (`j`). To swap the cards we use a tuple assignment:

```
1 (self.cards[i], self.cards[j]) = (self.cards[j], self.cards[i])
```

While this is a good shuffling method, a random number generator object also has a `shuffle` method that can shuffle elements in a list, in place. So we could rewrite this function to use the one provided for us:

```
1 class Deck:
2     ...
3     def shuffle(self):
4         import random
5         rng = random.Random()           # Create a random generator
6         rng.shuffle(self.cards)         # uUse its shuffle method
```

### 11.4.8 Removing and dealing cards

Another method that would be useful for the `Deck` class is `remove`, which takes a card as a parameter, removes it, and returns `True` if the card was in the deck and `False` otherwise:

```
1 class Deck:
2     ...
3     def remove(self, card):
4         if card in self.cards:
5             self.cards.remove(card)
6             return True
7         else:
8             return False
```

The `in` operator returns `True` if the first operand is in the second. If the first operand is an object, Python uses the object's `__eq__` method to determine equality with items in the list. Since the `__eq__` we provided in the `Card` class checks for deep equality, the `remove` method checks for deep equality.

To deal cards, we want to remove and return the top card. The list method `pop` provides a convenient way to do that:

```
1 class Deck:
2     ...
3     def pop(self):
4         return self.cards.pop()
```

Actually, `pop` removes the *last* card in the list, so we are in effect dealing from the bottom of the deck.

One more operation that we are likely to want is the Boolean function `is_empty`, which returns `True` if the deck contains no cards:

```
1 class Deck:
2     ...
3     def is_empty(self):
4         return self.cards == []
```

## 11.4.9 Glossary

**encode** To represent one type of value using another type of value by constructing a mapping between them.

**class attribute** A variable that is defined inside a class definition but outside any method. Class attributes are accessible from any method in the class and are shared by all instances of the class.

**accumulator** A variable used in a loop to accumulate a series of values, such as by concatenating them onto a string or adding them to a running sum.

## 11.4.10 Exercises

1. Modify `cmp` so that Aces are ranked higher than Kings.

## 11.5 Inheritance

### 11.5.1 Inheritance

The language feature most often associated with object-oriented programming is **inheritance**. Inheritance is the ability to define a new class that is a modified version of an existing class.

The primary advantage of this feature is that you can add new methods to a class without modifying the existing class. It is called inheritance because the new class inherits all of the methods of the existing class. Extending this metaphor, the existing class is sometimes called the **parent** class. The new class may be called the **child** class or sometimes subclass.

Inheritance is a powerful feature. Some programs that would be complicated without inheritance can be written concisely and simply with it. Also, inheritance can facilitate code reuse, since you can customize the behavior of parent classes without having to modify them. In some cases, the inheritance structure reflects the natural structure of the problem, which makes the program easier to understand.

On the other hand, inheritance can make programs difficult to read. When a method is invoked, it is sometimes not clear where to find its definition. The relevant code may be scattered among several modules. Also, many of the things that can be done using inheritance can be done as elegantly (or more so) without it. If the natural structure of the problem does not lend itself to inheritance, this style of programming can do more harm than good.

In this chapter we will demonstrate the use of inheritance as part of a program that plays the card game Old Maid. One of our goals is to write code that could be reused to implement other card games.

### 11.5.2 A hand of cards

For almost any card game, we need to represent a hand of cards. A hand is similar to a deck, of course. Both are made up of a set of cards, and both require operations like adding and removing cards. Also, we might like the ability to shuffle both decks and hands.

A hand is also different from a deck. Depending on the game being played, we might want to perform some operations on hands that don't make sense for a deck. For example, in poker we might classify a hand (straight, flush, etc.) or compare it with another hand. In bridge, we might want to compute a score for a hand in order to make a bid.

This situation suggests the use of inheritance. If `Hand` is a subclass of `Deck`, it will have all the methods of `Deck`, and new methods can be added.

We add the code in this chapter to our `Cards.py` file from the previous chapter. In the class definition, the name of the parent class appears in parentheses:

```
1 class Hand(Deck):
2     pass
```

This statement indicates that the new `Hand` class inherits from the existing `Deck` class.

The `Hand` constructor initializes the attributes for the hand, which are `name` and `cards`. The string `name` identifies this hand, probably by the name of the player that holds it. The name is an optional parameter with the empty string as a default value. `cards` is the list of cards in the hand, initialized to the empty list:

```
1 class Hand(Deck):
2     def __init__(self, name=""):
3         self.cards = []
4         self.name = name
```

For just about any card game, it is necessary to add and remove cards from the deck. Removing cards is already taken care of, since `Hand` inherits `remove` from `Deck`. But we have to write `add`:

```
1 class Hand(Deck):
2     ...
3     def add(self, card):
4         self.cards.append(card)
```

Again, the ellipsis indicates that we have omitted other methods. The list `append` method adds the new card to the end of the list of cards.

### 11.5.3 Dealing cards

Now that we have a `Hand` class, we want to deal cards from the `Deck` into hands. It is not immediately obvious whether this method should go in the `Hand` class or in the `Deck` class, but since it operates on a single deck and (possibly) several hands, it is more natural to put it in `Deck`.

`deal` should be fairly general, since different games will have different requirements. We may want to deal out the entire deck at once or add one card to each hand.

`deal` takes two parameters, a list (or tuple) of hands and the total number of cards to deal. If there are not enough cards in the deck, the method deals out all of the cards and stops:



```
1 class Deck:
2     ...
3     def deal(self, hands, num_cards=999):
4         num_hands = len(hands)
5         for i in range(num_cards):
6             if self.is_empty():
7                 break                # Break if out of cards
8             card = self.pop()         # Take the top card
9             hand = hands[i % num_hands] # Whose turn is next?
10            hand.add(card)            # Add the card to the hand
```

The second parameter, `num_cards`, is optional; the default is a large number, which effectively means that all of the cards in the deck will get dealt.

The loop variable `i` goes from 0 to `num_cards-1`. Each time through the loop, a card is removed from the deck using the list method `pop`, which removes and returns the last item in the list.

The modulus operator (`%`) allows us to deal cards in a round robin (one card at a time to each hand). When `i` is equal to the number of hands in the list, the expression `i % num_hands` wraps around to the beginning of the list (index 0).

### 11.5.4 Printing a Hand

To print the contents of a hand, we can take advantage of the `__str__` method inherited from `Deck`. For example:

```
>>> deck = Deck()
>>> deck.shuffle()
>>> hand = Hand("frank")
>>> deck.deal([hand], 5)
>>> print(hand)
Hand frank contains
2 of Spades
3 of Spades
4 of Spades
Ace of Hearts
9 of Clubs
```

It's not a great hand, but it has the makings of a straight flush.

Although it is convenient to inherit the existing methods, there is additional information in a `Hand` object we might want to include when we print one. To do that, we can provide a `__str__` method in the `Hand` class that overrides the one in the `Deck` class:

```
1 class Hand(Deck)
2     ...
3     def __str__(self):
4         s = "Hand " + self.name
5         if self.is_empty():
6             s += " is empty\n"
7         else:
8             s += " contains\n"
9         return s + Deck.__str__(self)
```

Initially, `s` is a string that identifies the hand. If the hand is empty, the program appends the words `is empty` and returns `s`.

Otherwise, the program appends the word `contains` and the string representation of the `Deck`, computed by invoking the `__str__` method in the `Deck` class on `self`.

It may seem odd to send `self`, which refers to the current `Hand`, to a `Deck` method, until you remember that a `Hand` is a kind of `Deck`. `Hand` objects can do everything `Deck` objects can, so it is legal to send a `Hand` to a `Deck` method.

In general, it is always legal to use an instance of a subclass in place of an instance of a parent class.

### 11.5.5 The `CardGame` class

The `CardGame` class takes care of some basic chores common to all games, such as creating the deck and shuffling it:

```
1 class CardGame:
2     def __init__(self):
3         self.deck = Deck()
4         self.deck.shuffle()
```

This is the first case we have seen where the initialization method performs a significant computation, beyond initializing attributes.

To implement specific games, we can inherit from `CardGame` and add features for the new game. As an example, we'll write a simulation of Old Maid.

The object of Old Maid is to get rid of cards in your hand. You do this by matching cards by rank and color. For example, the 4 of Clubs matches the 4 of Spades since both suits are black. The Jack of Hearts matches the Jack of Diamonds since both are red.

To begin the game, the Queen of Clubs is removed from the deck so that the Queen of Spades has no match. The fifty-one remaining cards are dealt to the players in a round robin. After the deal, all players match and discard as many cards as possible.

When no more matches can be made, play begins. In turn, each player picks a card (without looking) from the closest neighbor to the left who still has cards. If the chosen card matches a card in the player's hand, the pair is removed. Otherwise, the card is added to the player's hand. Eventually all possible matches are made, leaving only the Queen of Spades in the loser's hand.

In our computer simulation of the game, the computer plays all hands. Unfortunately, some nuances of the real game are lost. In a real game, the player with the Old Maid goes to some effort to get their neighbor to pick that card, by displaying it a little more prominently, or perhaps failing to display it more prominently, or even failing to fail to display that card more prominently. The computer simply picks a neighbor's card at random.

### 11.5.6 `OldMaidHand` class

A hand for playing Old Maid requires some abilities beyond the general abilities of a `Hand`. We will define a new class, `OldMaidHand`, that inherits from `Hand` and provides an additional method called `remove_matches`:

```
1 class OldMaidHand(Hand):
2     def remove_matches(self):
3         count = 0
4         original_cards = self.cards[:]
5         for card in original_cards:
6             match = Card(3 - card.suit, card.rank)
7             if match in self.cards:
8                 self.cards.remove(card)
9                 self.cards.remove(match)
10                print("Hand {0}: {1} matches {2}"
11                      .format(self.name, card, match))
12                count += 1
13        return count
```

We start by making a copy of the list of cards, so that we can traverse the copy while removing cards from the original. Since `self.cards` is modified in the loop, we don't want to use it to control the traversal. Python can get quite confused if it is traversing a list that is changing!

For each card in the hand, we figure out what the matching card is and go looking for it. The match card has the same rank and the other suit of the same color. The expression `3 - card.suit` turns a Club (suit 0) into a Spade (suit 3) and a Diamond (suit 1) into a Heart (suit 2). You should satisfy yourself that the opposite operations also work. If the match card is also in the hand, both cards are removed.

The following example demonstrates how to use `remove_matches`:

```
>>> game = CardGame()
>>> hand = OldMaidHand("frank")
>>> game.deck.deal([hand], 13)
>>> print(hand)
Hand frank contains
Ace of Spades
2 of Diamonds
7 of Spades
8 of Clubs
6 of Hearts
8 of Spades
7 of Clubs
Queen of Clubs
7 of Diamonds
5 of Clubs
Jack of Diamonds
10 of Diamonds
10 of Hearts
>>> hand.remove_matches()
Hand frank: 7 of Spades matches 7 of Clubs
Hand frank: 8 of Spades matches 8 of Clubs
Hand frank: 10 of Diamonds matches 10 of Hearts
>>> print(hand)
Hand frank contains
Ace of Spades
2 of Diamonds
6 of Hearts
Queen of Clubs
7 of Diamonds
5 of Clubs
Jack of Diamonds
```

Notice that there is no `__init__` method for the `OldMaidHand` class. We inherit it from `Hand`.

### 11.5.7 OldMaidGame class

Now we can turn our attention to the game itself. `OldMaidGame` is a subclass of `CardGame` with a new method called `play` that takes a list of players as a parameter.

Since `__init__` is inherited from `CardGame`, a new `OldMaidGame` object contains a new shuffled deck:

```
1 class OldMaidGame(CardGame):
2     def play(self, names):
3         # Remove Queen of Clubs
4         self.deck.remove(Card(0,12))
5
```

(continues on next page)

(continued from previous page)

```

6         # Make a hand for each player
7         self.hands = []
8         for name in names:
9             self.hands.append(OldMaidHand(name))
10
11        # Deal the cards
12        self.deck.deal(self.hands)
13        print("----- Cards have been dealt")
14        self.print_hands()
15
16        # Remove initial matches
17        matches = self.remove_all_matches()
18        print("----- Matches discarded, play begins")
19        self.print_hands()
20
21        # Play until all 50 cards are matched
22        turn = 0
23        num_hands = len(self.hands)
24        while matches < 25:
25            matches += self.play_one_turn(turn)
26            turn = (turn + 1) % num_hands
27
28        print("----- Game is Over")
29        self.print_hands()

```

The writing of `print_hands` has been left as an exercise.

Some of the steps of the game have been separated into methods. `remove_all_matches` traverses the list of hands and invokes `remove_matches` on each:

```

1 class OldMaidGame(CardGame):
2     ...
3     def remove_all_matches(self):
4         count = 0
5         for hand in self.hands:
6             count += hand.remove_matches()
7         return count

```

`count` is an accumulator that adds up the number of matches in each hand. When we've gone through every hand, the total is returned (`count`).

When the total number of matches reaches twenty-five, fifty cards have been removed from the hands, which means that only one card is left and the game is over.

The variable `turn` keeps track of which player's turn it is. It starts at 0 and increases by one each time; when it reaches `num_hands`, the modulus operator wraps it back around to 0.

The method `play_one_turn` takes a parameter that indicates whose turn it is. The return value is the number of matches made during this turn:

```

1 class OldMaidGame(CardGame):
2     ...
3     def play_one_turn(self, i):
4         if self.hands[i].is_empty():
5             return 0
6         neighbor = self.find_neighbor(i)
7         picked_card = self.hands[neighbor].pop()
8         self.hands[i].add(picked_card)

```

(continues on next page)

(continued from previous page)

```
9         print("Hand", self.hands[i].name, "picked", picked_card)
10        count = self.hands[i].remove_matches()
11        self.hands[i].shuffle()
12        return count
```

If a player's hand is empty, that player is out of the game, so he or she does nothing and returns 0.

Otherwise, a turn consists of finding the first player on the left that has cards, taking one card from the neighbor, and checking for matches. Before returning, the cards in the hand are shuffled so that the next player's choice is random.

The method `find_neighbor` starts with the player to the immediate left and continues around the circle until it finds a player that still has cards:

```
1 class OldMaidGame(CardGame):
2     ...
3     def find_neighbor(self, i):
4         num_hands = len(self.hands)
5         for next in range(1, num_hands):
6             neighbor = (i + next) % num_hands
7             if not self.hands[neighbor].is_empty():
8                 return neighbor
```

If `find_neighbor` ever went all the way around the circle without finding cards, it would return `None` and cause an error elsewhere in the program. Fortunately, we can prove that that will never happen (as long as the end of the game is detected correctly).

We have omitted the `print_hands` method. You can write that one yourself.

The following output is from a truncated form of the game where only the top fifteen cards (tens and higher) were dealt to three players. With this small deck, play stops after seven matches instead of twenty-five.

```
>>> import cards
>>> game = cards.OldMaidGame()
>>> game.play(["Allen", "Jeff", "Chris"])
----- Cards have been dealt
Hand Allen contains
King of Hearts
Jack of Clubs
Queen of Spades
King of Spades
10 of Diamonds

Hand Jeff contains
Queen of Hearts
Jack of Spades
Jack of Hearts
King of Diamonds
Queen of Diamonds

Hand Chris contains
Jack of Diamonds
King of Clubs
10 of Spades
10 of Hearts
10 of Clubs

Hand Jeff: Queen of Hearts matches Queen of Diamonds
Hand Chris: 10 of Spades matches 10 of Clubs
```

(continues on next page)

(continued from previous page)

```
----- Matches discarded, play begins
Hand Allen contains
King of Hearts
  Jack of Clubs
    Queen of Spades
      King of Spades
        10 of Diamonds

Hand Jeff contains
Jack of Spades
  Jack of Hearts
    King of Diamonds

Hand Chris contains
Jack of Diamonds
  King of Clubs
    10 of Hearts

Hand Allen picked King of Diamonds
Hand Allen: King of Hearts matches King of Diamonds
Hand Jeff picked 10 of Hearts
Hand Chris picked Jack of Clubs
Hand Allen picked Jack of Hearts
Hand Jeff picked Jack of Diamonds
Hand Chris picked Queen of Spades
Hand Allen picked Jack of Diamonds
Hand Allen: Jack of Hearts matches Jack of Diamonds
Hand Jeff picked King of Clubs
Hand Chris picked King of Spades
Hand Allen picked 10 of Hearts
Hand Allen: 10 of Diamonds matches 10 of Hearts
Hand Jeff picked Queen of Spades
Hand Chris picked Jack of Spades
Hand Chris: Jack of Clubs matches Jack of Spades
Hand Jeff picked King of Spades
Hand Jeff: King of Clubs matches King of Spades
----- Game is Over
Hand Allen is empty

Hand Jeff contains
Queen of Spades

Hand Chris is empty
```

So Jeff loses.

## 11.5.8 Glossary

**inheritance** The ability to define a new class that is a modified version of a previously defined class.

**parent class** The class from which a child class inherits.

**child class** A new class created by inheriting from an existing class; also called a subclass.

### 11.5.9 Exercises

1. Add a method, `print_hands`, to the `OldMaidGame` class which traverses `self.hands` and prints each hand.
2. Define a new kind of Turtle, `TurtleGTX`, that comes with some extra features: it can jump forward a given distance, and it has an odometer that keeps track of how far the turtle has travelled since it came off the production line. (The parent class has a number of synonyms like `fd`, `forward`, `back`, `backward`, and `bk`: for this exercise, just focus on putting this functionality into the `forward` method.) Think carefully about how to count the distance if the turtle is asked to move forward by a negative amount. (We would not want to buy a second-hand turtle whose odometer reading was faked because its previous owner drove it backwards around the block too often. Try this in a car near you, and see if the car's odometer counts up or down when you reverse.)
3. After travelling some random distance, your turtle should break down with a flat tyre. After this happens, raise an exception whenever `forward` is called. Also provide a `change_tyre` method that can fix the flat.

## 12.1 Catching exceptions

Whenever a runtime error occurs, it creates an **exception** object. The program stops running at this point and Python prints out the traceback, which ends with a message describing the exception that occurred.

For example, dividing by zero creates an exception:

```
>>> print(55/0)
Traceback (most recent call last):
  File "<interactive input>", line 1, in <module>
ZeroDivisionError: integer division or modulo by zero
```

So does accessing a non-existent list item:

```
>>> a = []
>>> print(a[5])
Traceback (most recent call last):
  File "<interactive input>", line 1, in <module>
IndexError: list index out of range
```

Or trying to make an item assignment on a tuple:

```
>>> tup = ("a", "b", "d", "d")
>>> tup[2] = "c"
Traceback (most recent call last):
  File "<interactive input>", line 1, in <module>
TypeError: 'tuple' object does not support item assignment
```

In each case, the error message on the last line has two parts: the type of error before the colon, and specifics about the error after the colon.

Sometimes we want to execute an operation that might cause an exception, but we don't want the program to stop. We can **handle the exception** using the `try` statement to “wrap” a region of code.



For example, we might prompt the user for the name of a file and then try to open it. If the file doesn't exist, we don't want the program to crash; we want to handle the exception:

```
1 filename = input("Enter a file name: ")
2 try:
3     f = open(filename, "r")
4 except FileNotFoundError:
5     print("There is no file named", filename)
```

The `try` statement has four separate clauses—or parts—introduced by the keywords `try`, `except`, `else`, and `finally`. All clauses but the `try` can be omitted.

The interpreter executes the block under the `try` statement, and monitors for exceptions. If one occurs, the interpreter moves to the `except` statement; it executes the `except` block if the exception raised match the exception requested in the `except` statement. If no exception occurs, the interpreter skips the block under the `except` clause. A `else` block is executed after the `try` one, if no exception occurred. A `finally` block is executed in any case. With all the statements, a `try` clause looks like:

```
1 user_input = input('Type a number:')
2 try:
3     # Try do do something that could fail.
4     user_input_as_number = float(user_input)
5 except ValueError:
6     # This will be executed if a ``ValueError`` is raised.
7     print('You did not enter a number.')
8 else:
9     # This will be executed if not exception got raised in the
10    # ``try`` statement.
11    print('The square of your number is ', user_input_as_number**2)
12 finally:
13    # This will be executed whether or not an exception is raised.
14    print('Thank you')
```

When using a `try` clause, you should have as little as possible in the `try` block. If too many things happen in that block, you risk handling an unexpected exception.

If the `try` block can fail if various way, you can handle different exceptions in the same `try` clause:

It is also possible not to specify a particular exception in the `except` statement. In this case, any exception will be handled. Such bare `except` statement should be avoided, though, as they can easily mask bugs.

## 12.2 Raising our own exceptions

Can our program deliberately cause its own exceptions? If our program detects an error condition, we can **raise** an exception. Here is an example that gets input from the user and checks that the number is non-negative:

```
1 def get_age():
2     age = int(input("Please enter your age: "))
3     if age < 0:
4         # Create a new instance of an exception
5         my_error = ValueError("{0} is not a valid age".format(age))
6         raise my_error
7     return age
```

Line 5 creates an exception object, in this case, a `ValueError` object, which encapsulates specific information about the error. Assume that in this case function A called B which called C which called D which called `get_age`. The `raise` statement on line 6 carries this object out as a kind of “return value”, and immediately exits from `get_age()` to its caller D. Then D again exits to its caller C, and C exits to B and so on, each returning the exception object to their caller, until it encounters a `try ... except` that can handle the exception. We call this “unwinding the call stack”.

`ValueError` is one of the built-in exception types which most closely matches the kind of error we want to raise. The complete listing of built-in exceptions can be found at the [Built-in Exceptions](#) section of the [Python Library Reference](#), again by Python’s creator, Guido van Rossum.

If the function that called `get_age` (or its caller, or their caller, ...) handles the error, then the program can carry on running; otherwise, Python prints the traceback and exits:

```
>>> get_age()
Please enter your age: 42
42
>>> get_age()
Please enter your age: -2
Traceback (most recent call last):
  File "<interactive input>", line 1, in <module>
  File "learn_exceptions.py", line 4, in get_age
    raise ValueError("{0} is not a valid age".format(age))
ValueError: -2 is not a valid age
```

The error message includes the exception type and the additional information that was provided when the exception object was first created.

It is often the case that lines 5 and 6 (creating the exception object, then raising the exception) are combined into a single statement, but there are really two different and independent things happening, so perhaps it makes sense to keep the two steps separate when we first learn to work with exceptions. Here we show it all in a single statement:

```
1 raise ValueError("{0} is not a valid age".format(age))
```

## 12.3 Revisiting an earlier example

Using exception handling, we can now modify our `recursion_depth` example from the previous chapter so that it stops when it reaches the maximum recursion depth allowed:

```
1 def recursion_depth(number):
2     print("Recursion depth number", number)
3     try:
4         recursion_depth(number + 1)
5     except:
6         print("I cannot go any deeper into this wormhole.")
7
8 recursion_depth(0)
```

Run this version and observe the results.

## 12.4 The `finally` clause of the `try` statement

A common programming pattern is to grab a resource of some kind — e.g. we create a window for turtles to draw on, or we dial up a connection to our internet service provider, or we may open a file for writing. Then we perform some

computation which may raise an exception, or may work without any problems.

Whatever happens, we want to “clean up” the resources we grabbed — e.g. close the window, disconnect our dial-up connection, or close the file. The `finally` clause of the `try` statement is the way to do just this. Consider this (somewhat contrived) example:

```
1 import turtle
2 import time
3
4 def show_poly():
5     try:
6         win = turtle.Screen()    # Grab/create a resource, e.g. a window
7         tess = turtle.Turtle()
8
9         # This dialog could be cancelled,
10        # or the conversion to int might fail, or n might be zero.
11        n = int(input("How many sides do you want in your polygon?"))
12        angle = 360 / n
13        for i in range(n):        # Draw the polygon
14            tess.forward(10)
15            tess.left(angle)
16        time.sleep(3)              # Make program wait a few seconds
17    finally:
18        win.bye()                  # Close the turtle's window
19
20 show_poly()
21 show_poly()
22 show_poly()
```

In lines 20–22, `show_poly` is called three times. Each one creates a new window for its turtle, and draws a polygon with the number of sides input by the user. But what if the user enters a string that cannot be converted to an `int`? What if they close the dialog? We’ll get an exception, *but even though we’ve had an exception, we still want to close the turtle’s window*. Lines 17–18 does this for us. Whether we complete the statements in the `try` clause successfully or not, the `finally` block will always be executed.

Notice that the exception is still unhandled — only an `except` clause can handle an exception, so our program will still crash. But at least its turtle window will be closed before it crashes!

## 12.5 Glossary

**exception** An error that occurs at runtime.

**handle an exception** To prevent an exception from causing our program to crash, by wrapping the block of code in a `try... except` construct.

**raise** To create a deliberate exception by using the `raise` statement.

## 12.6 Exercises

1. Write a function named `readposint` that uses the `input` dialog to prompt the user for a positive integer and then checks the input to confirm that it meets the requirements. It should be able to handle inputs that cannot be converted to `int`, as well as negative `ints`, and edge cases (e.g. when the user closes the dialog, or does not enter anything at all.)

## CHAPTER 13

---

### Fitting

---

Suppose we want to determine the gravitational acceleration. To this end, we could drop an object from the building and measure how long it takes for the object to reach the ground with a stopwatch. Newton's laws predict the following model:

$$h = \frac{1}{2}gt^2$$

where  $h$  is the height from which we dropped the object, and  $t$  the time it takes to hit the ground. So from our one measurement we could now calculate the gravitational acceleration  $g$ .

We can measure the height very accurately, but since we use a stopwatch to measure time, that value is a lot less reliable because you might have started and stopped your stopwatch at the wrong moments. Therefore, the result will not be very accurate. To make the value more accurate we should repeat the same measurement  $n$  times to obtain an average  $t$  and use that instead. We'll get back to this later.

For now we are more interested in the question: is this model correct? To test this question, we drop our object from different heights, doing multiple measurements for each height to get reliable values. The data, obtained by simulation for health and safety reasons, are given in the following table:

y	t	n
10	1.4	5
20	2.1	3
30	2.6	8
40	3.0	15
50	3.3	30

Since the model predicts a parabola, we want to fit the data to this model to see how good it works. It might be a bit confusing, but  $h$  is our x axis, and  $t$  is the y axis.

We use the `symfit` package to do our fitting. You can find the installation instructions [here](#).

To fit the data to the model we run the following code:

```
import numpy as np
from sympfit import Variable, Parameter, Fit, Model, sqrt

t_data = np.array([1.4, 2.1, 2.6, 3.0, 3.3])
h_data = np.array([10, 20, 30, 40, 50])

# We now define our model
h = Variable('h')
t = Variable('t')
g = Parameter('g')

t_model = Model({t: sqrt(2 * h / g)})

fit = Fit(t_model, h=h_data, t=t_data)
fit_result = fit.execute()
print(fit_result)
```

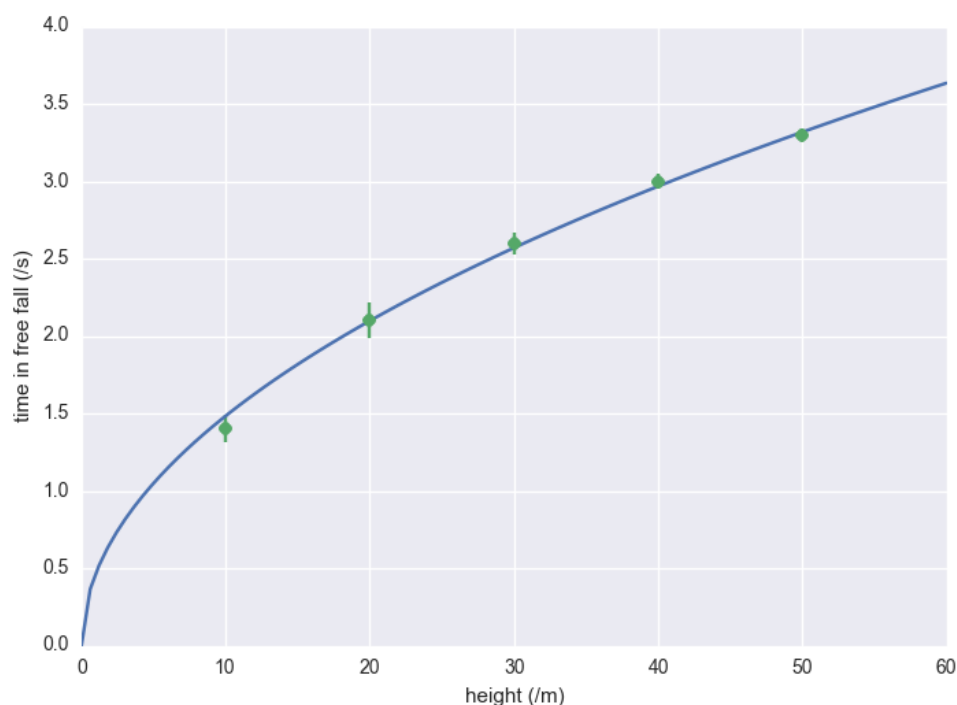
Looking at these results, we see that  $g = 9.09 \pm 0.15$  for this dataset. In order to plot this result alongside the data, we need to calculate values for the model. In the same script, we can do:

```
# Make an array from 0 to 50 in 1000 steps
h_range = np.linspace(0, 50, 1000)
fit_data = t_model(h=h_range, g=fit_result.value(g))
t_fit = fit_data.t
```

---

**Note:** When calling a `sympfit.Model`, a `namedtuple` is returned with all the components of the model evaluated. In this case there is only one component, and can be accessed by requesting `fit_data.t` or `fit_data[0]`.

---



This gives the model evaluated at all the points in `h_range`. Making the actual plot is left to you as an exercise. We see that we can request the fitted value of `g` by calling `fit_result.value(g)`, which returns 9.09.

Let's think for a second about the implications. The value of  $g$  is  $g = 9.81$  in the Netherlands. Based on our result, the textbooks should be rewritten because that value is extremely unlikely to be true given the small standard deviation in our data. It is at this point that we remember that our data point were not infinitely precise: we took many measurements and averaged them. This means there is an uncertainty in each of our data points. We will now account for this additional uncertainty and see what this does to our conclusion. To do this we first have to describe how the fitting actually works.

## 13.1 How does it work?

In fitting we want to find values for the parameters such that the differences between the model and the data are as small as possible. The differences (residuals) are easy to calculate:

$$f(x_i, \vec{p}) - y_i$$

Here we have written the parameters as a vector  $\vec{p}$ , to indicate that we can have multiple parameters.  $x_i$  and  $y_i$  are the  $x$  and  $y$  coordinates of the  $i$ 'th datapoint. However, if we were to minimize the sum over all these differences we would have a problem, because these differences can be either positive or negative. This means there's many ways to add these values and get zero out of the sum. We therefore take the sum over the residuals squared:

$$Q^2 = \sum_{i=1}^n (f(x_i, \vec{p}) - y_i)^2$$

Now if we minimize  $Q^2$ , we get the best possible values for our parameters. The fitting algorithm actually just takes some values for the parameters, calculates  $Q^2$ , then changes the values slightly by adding or subtracting a number, and checks if this new value is smaller than the old one. If this is true it keeps going in the same direction until the value of  $Q^2$  starts to increase. That's when you know you've hit a minimum. Of course the trick is to do this smartly, and a lot of algorithms have been developed in order to do this.

## 13.2 Propagating Uncertainties

In the example above we the fitting process assumed that every measurement was equally reliable. But this is not true. By repeating a measurement and averaging the result, we can improve the accuracy. So in our example, we dropped our object from every height a couple of times and took the average. Therefore, we want to assign a weight depending on how accurate the average value for that height is. Statistically the weight  $w_i$  to use is  $w_i = \frac{1}{\sigma_i^2}$ , where  $\sigma_i$  is the standard deviation for each point.

Our sum to minimize now changes to:

$$\chi^2 = \sum_{i=1}^n w_i (f(x_i, \vec{p}) - y_i)^2 = \sum_{i=1}^n \frac{(f(x_i, \vec{p}) - y_i)^2}{\sigma_i^2}$$

But how do we know the standard deviation in the mean value we calculate for every height? Suppose the standard deviation of our stopwatch is  $\sigma_{stopwatch} = 0.2$ . If we do  $n$  measurements from the same height, the average time is found by calculating

$$\bar{t} = \frac{1}{n} \sum_{i=1}^n t_i$$

It can be shown that the standard deviation of the mean is now:

$$\sigma_{\bar{t}} = \frac{\sigma_{stopwatch}}{\sqrt{n}}$$

So we see that by increasing the amount of measurements, we can decrease the uncertainty in  $\bar{t}$ . Our simulated data now changes to:

y	t	n	$\sigma_t$
10	1.4	5	0.089
20	2.1	3	0.115
30	2.6	8	0.071
40	3.0	15	0.052
50	3.3	30	0.037

The values of  $\sigma_t$  have been calculated by using the above formula. Let's fit to this new data set using *symfit*. Notice that there are some small differences to the code:

```
import numpy as np
from symfit import Variable, Parameter, Fit, Model, sqrt

t_data = np.array([1.4, 2.1, 2.6, 3.0, 3.3])
h_data = np.array([10, 20, 30, 40, 50])
n = np.array([5, 3, 8, 15, 30])
sigma = 0.2
sigma_t = sigma / np.sqrt(n)

# We now define our model
h = Variable('h')
t = Variable('t')
g = Parameter('g')
t_model = Model({t: sqrt(2 * h / g)})

fit = Fit(t_model, h=h_data, t=t_data, sigma_t=sigma_t)
fit_result = fit.execute()
print(fit_result)
```

---

**Note:** Looking at the initiation of `Fit`, we see that uncertainties can be provided to `Variable`'s by prepending their name with `sigma_`, in this case `sigma_t`.

---

Including these uncertainties in the fit yields  $g = 9.10 \pm 0.16$ . The accepted value of  $g = 9.81$  is well outside the uncertainty in this data. Therefore the textbooks must be rewritten!

This example shows the importance of propagating your errors consistently. (And of the importance of performing the actual measurement as the author of a chapter on error propagation so you don't end up claiming the textbooks have to be rewritten.)

## 13.3 More on symfit

There are a lot more features in `symfit` to help you on your quest to fitting the universe. You can find the tutorial [here](#).

It is recommended you read this as well before starting to fit your own data.



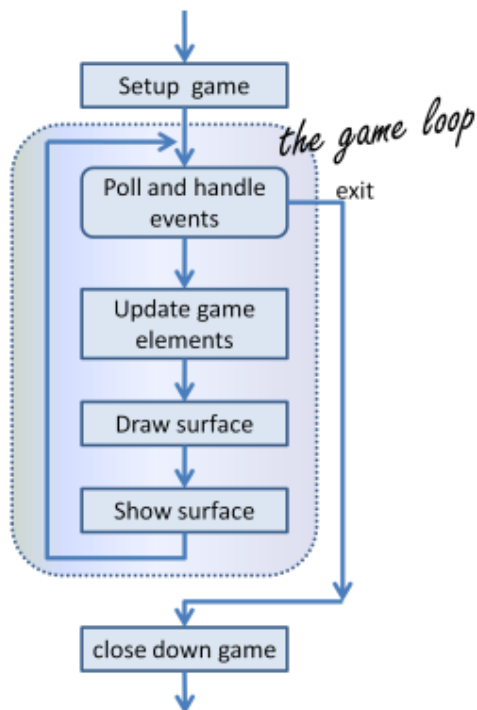


PyGame is a package that is not part of the standard Python distribution, so if you do not already have it installed (i.e. `import pygame` fails), download and install a suitable version from <http://pygame.org/download.shtml>. These notes are based on PyGame 1.9.1, the most recent version at the time of writing.

PyGame comes with a substantial set of tutorials, examples, and help, so there is ample opportunity to stretch yourself on the code. You may need to look around a bit to find these resources, though: if you've installed PyGame on a Windows machine, for example, they'll end up in a folder like `C:\Python31\Lib\site-packages\pygame\` where you will find directories for *docs* and *examples*.

## 14.1 The game loop

The structure of the games we'll consider always follows this fixed pattern:



In every game, in the *setup* section we'll create a window, load and prepare some content, and then enter the **game loop**. The game loop continuously does four main things:

- it **polls** for events — i.e. asks the system whether events have occurred — and responds appropriately,
- it updates whatever internal data structures or objects need changing,
- it draws the current state of the game into a (non-visible) surface,
- it puts the just-drawn surface on display.

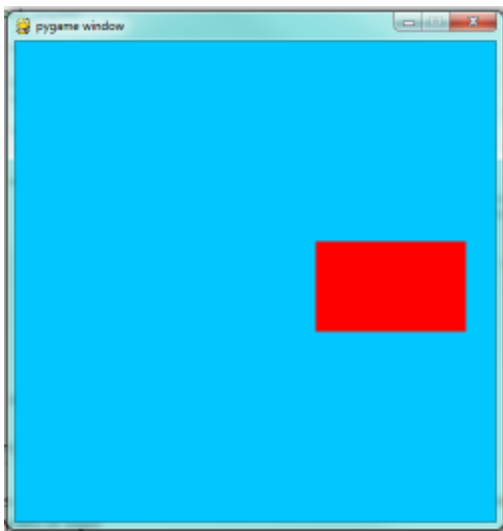
```
1  import pygame
2
3  def main():
4      """ Set up the game and run the main game loop """
5      pygame.init()          # Prepare the pygame module for use
6      surface_size = 480     # Desired physical surface size, in pixels.
7
8      # Create surface of (width, height), and its window.
9      main_surface = pygame.display.set_mode((surface_size, surface_size))
10
11     # Set up some data to describe a small rectangle and its color
12     small_rect = (300, 200, 150, 90)
13     some_color = (255, 0, 0)    # A color is a mix of (Red, Green, Blue)
14
15     while True:
16         event = pygame.event.poll()    # Look for any event
17         if event.type == pygame.QUIT:  # Window close button clicked?
18             break                      # ... leave game loop
19
20         # Update your game objects and data structures here...
21
22         # We draw everything from scratch on each frame.
23         # So first fill everything with the background color
```

(continues on next page)

(continued from previous page)

```
24     main_surface.fill((0, 200, 255))
25
26     # Overpaint a smaller rectangle on the main surface
27     main_surface.fill(some_color, small_rect)
28
29     # Now the surface is ready, tell pygame to display it!
30     pygame.display.flip()
31
32     pygame.quit()      # Once we leave the loop, close the window.
33
34 main()
```

This program pops up a window which stays there until we close it:



PyGame does all its drawing onto rectangular *surfaces*. After initializing PyGame at line 5, we create a window holding our main surface. The main loop of the game extends from line 15 to 30, with the following key bits of logic:

- First (line 16) we poll to fetch the next event that might be ready for us. This step will always be followed by some conditional statements that will determine whether any event that we’re interested in has happened. Polling for the event consumes it, as far as PyGame is concerned, so we only get one chance to fetch and use each event. On line 17 we test whether the type of the event is the predefined constant called `pygame.QUIT`. This is the event that we’ll see when the user clicks the close button on the PyGame window. In response to this event, we leave the loop.
- Once we’ve left the loop, the code at line 32 closes window, and we’ll return from function `main`. Your program could go on to do other things, or reinitialize pygame and create another window, but it will usually just end too.
- There are different kinds of events — key presses, mouse motion, mouse clicks, joystick movement, and so on. It is usual that we test and handle all these cases with new code squeezed in before line 19. The general idea is “handle events first, then worry about the other stuff”.
- At line 20 we’d update objects or data — for example, if we wanted to vary the color, position, or size of the rectangle we’re about to draw, we’d re-assign `some_color`, and `small_rect` here.
- A modern way to write games (now that we have fast computers and fast graphics cards) is to redraw everything from scratch on every iteration of the game loop. So the first thing we do at line 24 is fill the entire surface with a background color. The `fill` method of a surface takes two arguments — the color to use for filling, and the rectangle to be filled. But the second argument is optional, and if it is left out the entire surface is filled.

- In line 27 we fill a second rectangle, this time using `some_color`. The placement and size of the rectangle are given by the tuple `small_rect`, a 4-element tuple `(x, y, width, height)`.
- It is important to understand that the origin of the PyGame's surface is at the top left corner (unlike the turtle module that puts its origin in the middle of the screen). So, if you wanted the rectangle closer to the top of the window, you need to make its y coordinate smaller.
- If your graphics display hardware tries to read from memory at the same time as the program is writing to that memory, they will interfere with each other, causing video noise and flicker. To get around this, PyGame keeps two buffers in the main surface — the *back buffer* that the program draws to, while the *front buffer* is being shown to the user. Each time the program has fully prepared its back buffer, it flips the back/front role of the two buffers. So the drawing on lines 24 and 27 does not change what is seen on the screen until we `flip` the buffers, on line 30.

## 14.2 Displaying images and text

To draw an image on the main surface, we load the image, say a beach ball, into its own new surface. The main surface has a `blit` method that copies pixels from the beach ball surface into its own surface. When we call `blit`, we can specify where the beach ball should be placed on the main surface. The term **blit** is widely used in computer graphics, and means *to make a fast copy of pixels from one area of memory to another*.

So in the setup section, before we enter the game loop, we'd load the image, like this:

```
1 ball = pygame.image.load("ball.png")
```

and after line 28 in the program above, we'd add this code to display our image at position (100,120):

```
1 main_surface.blit(ball, (100, 120))
```

To display text, we need do do three things. Before we enter the game loop, we instantiate a `font` object:

```
1 # Instantiate 16 point Courier font to draw text.
2 my_font = pygame.font.SysFont("Courier", 16)
```

and after line 28, again, we use the font's `render` method to create a new surface containing the pixels of the drawn text, and then, as in the case for images, we `blit` our new surface onto the main surface. Notice that `render` takes two extra parameters — the second tells it whether to carefully smooth edges of the text while drawing (this process is called *anti-aliasing*), and the second is the color that we want the text to be. Here we've used `(0, 0, 0)` which is black:

```
1 the_text = my_font.render("Hello, world!", True, (0,0,0))
2 main_surface.blit(the_text, (10, 10))
```

We'll demonstrate these two new features by counting the frames — the iterations of the game loop — and keeping some timing information. On each frame, we'll display the frame count, and the frame rate. We will only update the frame rate after every 500 frames, when we'll look at the timing interval and can do the calculations.

```
1 import pygame
2 import time
3
4 def main():
5
6     pygame.init()      # Prepare the PyGame module for use
7     main_surface = pygame.display.set_mode((480, 240))
8
```

(continues on next page)

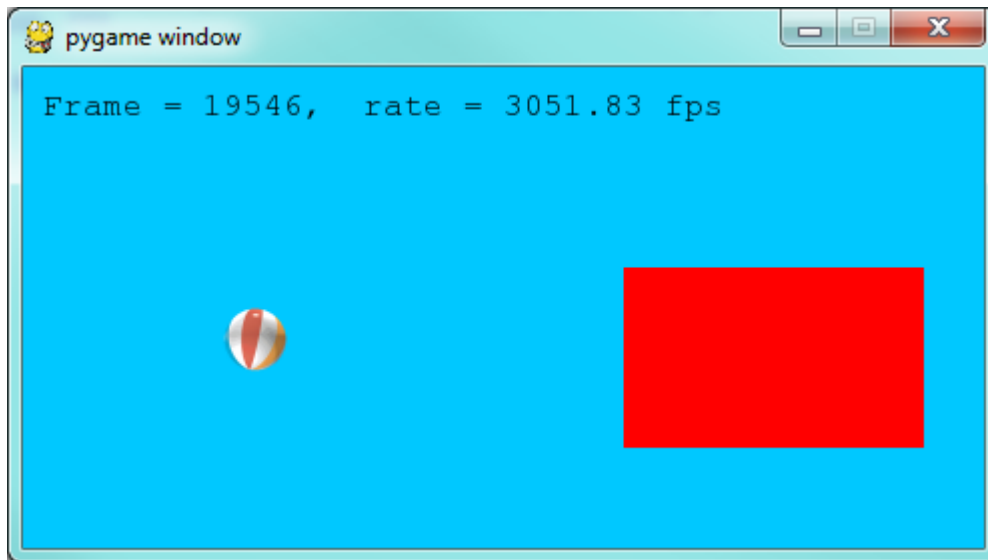
(continued from previous page)

```

9      # Load an image to draw. Substitute your own.
10     # PyGame handles gif, jpg, png, etc. image types.
11     ball = pygame.image.load("ball.png")
12
13     # Create a font for rendering text
14     my_font = pygame.font.SysFont("Courier", 16)
15
16     frame_count = 0
17     frame_rate = 0
18     t0 = time.clock()
19
20     while True:
21
22         # Look for an event from keyboard, mouse, joystick, etc.
23         ev = pygame.event.poll()
24         if ev.type == pygame.QUIT:    # Window close button clicked?
25             break                    # Leave game loop
26
27         # Do other bits of logic for the game here
28         frame_count += 1
29         if frame_count % 500 == 0:
30             t1 = time.clock()
31             frame_rate = 500 / (t1-t0)
32             t0 = t1
33
34         # Completely redraw the surface, starting with background
35         main_surface.fill((0, 200, 255))
36
37         # Put a red rectangle somewhere on the surface
38         main_surface.fill((255,0,0), (300, 100, 150, 90))
39
40         # Copy our image to the surface, at this (x,y) posn
41         main_surface.blit(ball, (100, 120))
42
43         # Make a new surface with an image of the text
44         the_text = my_font.render("Frame = {0}, rate = {1:.2f} fps"
45                                   .format(frame_count, frame_rate), True, (0,0,0))
46         # Copy the text surface to the main surface
47         main_surface.blit(the_text, (10, 10))
48
49         # Now that everything is drawn, put it on display!
50         pygame.display.flip()
51
52     pygame.quit()
53
54
55     main()

```

The frame rate is close to ridiculous — a lot faster than one’s eye can process frames. (Commercial video games usually plan their action for 60 frames per second (fps).) Of course, our rate will drop once we start doing something a little more strenuous inside our game loop.



## 14.3 Drawing a board for the N queens puzzle

We previously solved our N queens puzzle. For the 8x8 board, one of the solutions was the list `[6, 4, 2, 0, 5, 7, 1, 3]`. Let's use that solution as testdata, and now use PyGame to draw that chessboard with its queens.

We'll create a new module for the drawing code, called `draw_queens.py`. When we have our test case(s) working, we can go back to our solver, import this new module, and add a call to our new function to draw a board each time a solution is discovered.

We begin with a background of black and red squares for the board. Perhaps we could create an image that we could load and draw, but that approach would need different background images for different size boards. Just drawing our own red and black rectangles of the appropriate size sounds like much more fun!

```
1 def draw_board(the_board):
2     """ Draw a chess board with queens, from the_board. """
3
4     pygame.init()
5     colors = [(255,0,0), (0,0,0)]    # Set up colors [red, black]
6
7     n = len(the_board)                # This is an NxN chess board.
8     surface_size = 480                # Proposed physical surface size.
9     square_size = surface_size // n   # sq_sz is length of a square.
10    surface_size = n * square_size     # Adjust to exactly fit n squares.
11
12    # Create the surface of (width, height), and its window.
13    surface = pygame.display.set_mode((surface_size, surface_size))
```

Here we precompute `square_size`, the integer size that each square will be, so that we can fit the squares nicely into the available window. So if we'd like the board to be 480x480, and we're drawing an 8x8 chessboard, then each square will need to have a size of 60 units. But we notice that a 7x7 board cannot fit nicely into 480 — we're going to get some ugly border that our squares don't fill exactly. So we recompute the surface size to exactly fit our squares before we create the window.

Now let's draw the squares, in the game loop. We'll need a nested loop: the outer loop will run over the rows of the chessboard, the inner loop over the columns:

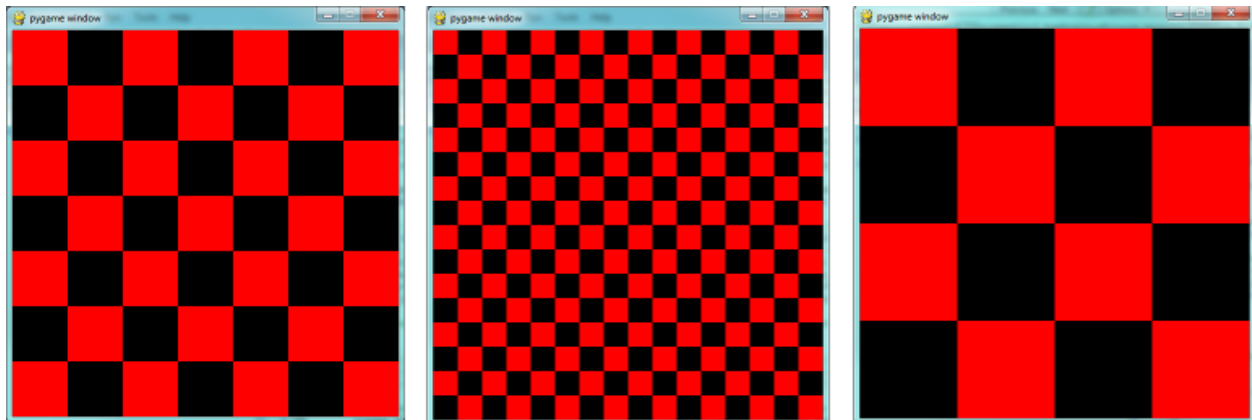
```

1  # Draw a fresh background (a blank chess board)
2  for row in range(n):          # Draw each row of the board.
3      color_index = row % 2      # Change starting color on each row
4      for col in range(n):      # Run through cols drawing squares
5          the_square = (col*square_size, row*square_size, square_size, square_
6      → size)
7          surface.fill(colors[color_index], the_square)
7          # now flip the color index for the next square
8          c_index = (c_index + 1) % 2

```

There are two important ideas in this code: firstly, we compute the rectangle to be filled from the `row` and `col` loop variables, multiplying them by the size of the square to get their position. And, of course, each square is a fixed width and height. So `the_square` represents the rectangle to be filled on the current iteration of the loop. The second idea is that we have to alternate colors on every square. In the earlier setup code we created a list containing two colors, here we manipulate `color_index` (which will always either have the value 0 or 1) to start each row on a color that is different from the previous row's starting color, and to switch colors each time a square is filled.

This (together with the other fragments not shown to flip the surface onto the display) leads to the pleasing backgrounds like this, for different size boards:



Now, on to drawing the queens! Recall that our solution `[6, 4, 2, 0, 5, 7, 1, 3]` means that in column 0 of the board we want a queen at row 6, at column 1 we want a queen at row 4, and so on. So we need a loop running over each queen:

```

1  for (col, row) in enumerate(the_board):
2      # draw a queen at col, row...

```

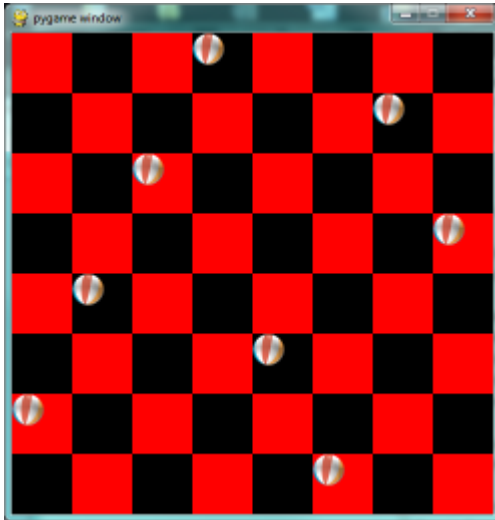
In this chapter we already have a beach ball image, so we'll use that for our queens. In the setup code before our game loop, we load the ball image (as we did before), and in the body of the loop, we add the line:

```

1  surface.blit(ball, (col * square_size, row * square_size))

```





We're getting there, but those queens need to be centred in their squares! Our problem arises from the fact that both the ball and the rectangle have their upper left corner as their reference points. If we're going to centre this ball in the square, we need to give it an extra offset in both the x and y direction. (Since the ball is round and the square is square, the offset in the two directions will be the same, so we'll just compute a single offset value, and use it in both directions.)

The offset we need is half the (size of the square less the size of the ball). So we'll precompute this in the game's setup section, after we've loaded the ball and determined the square size:

```
1 ball_offset = (square_size - ball.get_width()) // 2
```

Now we touch up the drawing code for the ball and we're done:

```
1 surface.blit(ball, (col * square_size + ball_offset, row * square_size +  
↪ball_offset))
```

We might just want to think about what would happen if the ball was bigger than the square. In that case, `ball_offset` would become negative. So it would still be centered in the square - it would just spill over the boundaries, or perhaps obscure the square entirely!

Here is the complete program:

```
1 import pygame
2
3 def draw_board(the_board):
4     """ Draw a chess board with queens, as determined by the the_board. """
5
6     pygame.init()
7     colors = [(255,0,0), (0,0,0)]    # Set up colors [red, black]
8
9     n = len(the_board)               # This is an NxN chess board.
10    surface_size = 480                # Proposed physical surface size.
11    square_size = surface_size // n   # sq_sz is length of a square.
12    surface_size = n * square_size    # Adjust to exactly fit n squares.
13
14    # Create the surface of (width, height), and its window.
15    surface = pygame.display.set_mode((surface_size, surface_size))
16
17    ball = pygame.image.load("ball.png")
18
```

(continues on next page)

(continued from previous page)

```

19     # Use an extra offset to centre the ball in its square.
20     # If the square is too small, offset becomes negative,
21     # but it will still be centered :-)
22     ball_offset = (square_size-ball.get_width()) // 2
23
24     while True:
25
26         # Look for an event from keyboard, mouse, etc.
27         event = pygame.event.poll()
28         if event.type == pygame.QUIT:
29             break;
30
31         # Draw a fresh background (a blank chess board)
32         for row in range(n):          # Draw each row of the board.
33             color_index = row % 2     # Alternate starting color
34             for col in range(n):      # Run through cols drawing squares
35                 the_square = (col*square_size, row*square_size, square_size,
↪square_size)
36                 surface.fill(colors[color_index], the_square)
37                 # Now flip the color index for the next square
38                 color_index = (color_index + 1) % 2
39
40         # Now that squares are drawn, draw the queens.
41         for (col, row) in enumerate(the_board):
42             surface.blit(ball,
43                           (col*square_size+ball_offset, row*square_size+ball_offset))
44
45         pygame.display.flip()
46
47
48     pygame.quit()
49
50 if __name__ == "__main__":
51     draw_board([0, 5, 3, 1, 6, 4, 2])    # 7 x 7 to test window size
52     draw_board([6, 4, 2, 0, 5, 7, 1, 3])
53     draw_board([9, 6, 0, 3, 10, 7, 2, 4, 12, 8, 11, 5, 1]) # 13 x 13
54     draw_board([11, 4, 8, 12, 2, 7, 3, 15, 0, 14, 10, 6, 13, 1, 5, 9])

```

There is one more thing worth reviewing here. The conditional statement on line 50 tests whether the name of the currently executing program is `__main__`. This allows us to distinguish whether this module is being run as a main program, or whether it has been imported elsewhere, and used as a module. If we run this module in Python, the test cases in lines 51-54 will be executed. However, if we import this module into another program (i.e. our N queens solver from earlier) the condition at line 50 will be false, and the statements on lines 51-54 won't run.

Previously, our main program looked like this:

```

1 def main():
2
3     board = list(range(8))    # Generate the initial permutation
4     num_found = 0
5     tries = 0
6     while num_found < 10:
7         random.shuffle(board)
8         tries += 1
9         if not has_clashes(board):
10             print("Found solution {0} in {1} tries.".format(board, tries))
11             tries = 0

```

(continues on next page)

(continued from previous page)

```
12         num_found += 1
13
14     main()
```

Now we just need two changes. At the top of that program, we import the module that we’ve been working on here (assume we called it `draw_queens`). (You’ll have to ensure that the two modules are saved in the same folder.) Then after line 10 here we add a call to draw the solution that we’ve just discovered:

```
draw_queens.draw_board(board)
```

And that gives a very satisfying combination of program that can search for solutions to the N queens problem, and when it finds each, it pops up the board showing the solution.

## 14.4 Sprites

A sprite is an object that can move about in a game, and has internal behaviour and state of its own. For example, a spaceship would be a sprite, the player would be a sprite, and bullets and bombs would all be sprites.

Object oriented programming (OOP) is ideally suited to a situation like this: each object can have its own attributes and internal state, and a couple of methods. Let’s have some fun with our N queens board. Instead of placing the queen in her final position, we’d like to drop her in from the top of the board, and let her fall into position, perhaps bouncing along the way.

The first encapsulation we need is to turn each of our queens into an object. We’ll keep a list of all the active sprites (i.e. a list of queen objects), and arrange two new things in our game loop:

- After handling events, but before drawing, call an `update` method on every sprite. This will give each sprite a chance to modify its internal state in some way — perhaps change its image, or change its position, or rotate itself, or make itself grow a bit bigger or a bit smaller.
- Once all the sprites have updated themselves, the game loop can begin drawing - first the background, and then call a `draw` method on each sprite in turn, and delegate (hand off) the task of drawing to the object itself. This is in line with the OOP idea that we don’t say “Hey, *draw*, show this queen!”, but we prefer to say “Hey, *queen*, draw yourself!”.

We start with a simple object, no movement or animation yet, just scaffolding, to see how to fit all the pieces together:

```
1 class QueenSprite:
2
3     def __init__(self, img, target_posn):
4         """ Create and initialize a queen for this
5             target position on the board
6         """
7         self.image = img
8         self.target_posn = target_posn
9         self.position = target_posn
10
11     def update(self):
12         return # Do nothing for the moment.
13
14     def draw(self, target_surface):
15         target_surface.blit(self.image, self.position)
```

We’ve given the sprite three attributes: an image to be drawn, a target position, and a current position. If we’re going to move the spite about, the current position may need to be different from the target, which is where we want the queen

finally to end up. In this code at this time we've done nothing in the `update` method, and our `draw` method (which can probably remain this simple in future) simply draws itself at its current position on the surface that is provided by the caller.

With its class definition in place, we now instantiate our N queens, put them into a list of sprites, and arrange for the game loop to call the `update` and `draw` methods on each frame. The new bits of code, and the revised game loop look like this:

```

1  all_sprites = []           # Keep a list of all sprites in the game
2
3  # Create a sprite object for each queen, and populate our list.
4  for (col, row) in enumerate(the_board):
5      a_queen = QueenSprite(ball,
6                             (col*square_size+ball_offset, row*square_size+ball_
7                             ↳offset))
8      all_sprites.append(a_queen)
9
10 while True:
11     # Look for an event from keyboard, mouse, etc.
12     event = pygame.event.poll()
13     if event.type == pygame.QUIT:
14         break;
15
16     # Ask every sprite to update itself.
17     for sprite in all_sprites:
18         sprite.update()
19
20     # Draw a fresh background (a blank chess board)
21     # ... same as before ...
22
23     # Ask every sprite to draw itself.
24     for sprite in all_sprites:
25         sprite.draw(surface)
26
27     pygame.display.flip()
```

This works just like it did before, but our extra work in making objects for the queens has prepared the way for some more ambitious extensions.

Let us begin with a falling queen object. At any instant, it will have a velocity i.e. a speed, in a certain direction. (We are only working with movement in the y direction, but use your imagination!) So in the object's `update` method, we want to change its current position by its velocity. If our N queens board is floating in space, velocity would stay constant, but hey, here on Earth we have gravity! Gravity changes the velocity on each time interval, so we'll want a ball that speeds up as it falls further. Gravity will be constant for all queens, so we won't keep it in the instances — we'll just make it a variable in our module. We'll make one other change too: we will start every queen at the top of the board, so that it can fall towards its target position. With these changes, we now get the following:

```

1  gravity = 0.0001
2
3  class QueenSprite:
4
5      def __init__(self, img, target_posn):
6          self.image = img
7          self.target_position = target_position
8          (x, y) = target_position
9          self.position = (x, 0)      # Start ball at top of its column
10         self.y_velocity = 0         # with zero initial velocity
11
```

(continues on next page)

(continued from previous page)

```
12 def update(self):
13     self.y_velocity += gravity          # Gravity changes velocity
14     (x, y) = self.position
15     new_y_pos = y + self.y_velocity    # Velocity moves the ball
16     self.position = (x, new_y_pos)     # to this new position.
17
18 def draw(self, target_surface):        # Same as before.
19     target_surface.blit(self.image, self.position)
```

Making these changes gives us a new chessboard in which each queen starts at the top of its column, and speeds up, until it drops off the bottom of the board and disappears forever. A good start — we have movement!

The next step is to get the ball to bounce when it reaches its own target position. It is pretty easy to bounce something — you just change the sign of its velocity, and it will move at the same speed in the opposite direction. Of course, if it is travelling up towards the top of the board it will be slowed down by gravity. (Gravity always sucks down!) And you'll find it bounces all the way up to where it began from, reaches zero velocity, and starts falling all over again. So we'll have bouncing balls that never settle.

A realistic way to settle the object is to lose some energy (probably to friction) each time it bounces — so instead of simply reversing the sign of the velocity, we multiply it by some fractional factor — say -0.65. This means the ball only retains 65% of its energy on each bounce, so it will, as in real life, stop bouncing after a short while, and settle on its “ground”.

The only changes are in the update method, which now looks like this:

```
1 def update(self):
2     self.y_velocity += gravity
3     (x, y) = self.position
4     new_y_pos = y + self.y_velocity
5     (target_x, target_y) = self.target_posn    # Unpack the position
6     dist_to_go = target_y - new_y_pos          # How far to our floor?
7
8     if dist_to_go < 0:                          # Are we under floor?
9         self.y_velocity = -0.65 * self.y_velocity    # Bounce
10        new_y_pos = target_y + dist_to_go          # Move back above floor
11
12    self.position = (x, new_y_pos)              # Set our new position.
```

Heh, heh, heh! We're not going to show animated screenshots, so copy the code into your Python environment and see for yourself.

## 14.5 Events

The only kind of event we're handled so far has been the QUIT event. But we can also detect keydown and keyup events, mouse motion, and mousebutton down or up events. Consult the PyGame documentation and follow the link to Event.

When your program polls for and receives an event object from PyGame, its event type will determine what secondary information is available. Each event object carries a *dictionary* (which you may only cover in due course in these notes). The dictionary holds certain *keys* and *values* that make sense for the type of event.

For example, if the type of event is MOUSEMOTION, we'll be able to find the mouse position and information about the state of the mouse buttons in the dictionary attached to the event. Similarly, if the event is KEYDOWN, we can learn from the dictionary which key went down, and whether any modifier keys (shift, control, alt, etc.) are also down. You also get events when the game window becomes active (i.e. gets focus) or loses focus.

The event object with type `NOEVENT` is returned if there are no events waiting. Events can be printed, allowing you to experiment and play around. So dropping these lines of code into the game loop directly after polling for any event is quite informative:

```
1  if event.type != pygame.NOEVENT:    # Only print if it is interesting!
2      print(event)
```

With this in place, hit the space bar and the escape key, and watch the events you get. Click your three mouse buttons. Move your mouse over the window. (This causes a vast cascade of events, so you may also need to filter those out of the printing.) You'll get output that looks something like this:

```
<Event(17-VideoExpose {})>
<Event(1-ActiveEvent {'state': 1, 'gain': 0})>
<Event(2-KeyDown {'scancode': 57, 'key': 32, 'unicode': ' ', 'mod': 0})>
<Event(3-KeyUp {'scancode': 57, 'key': 32, 'mod': 0})>
<Event(2-KeyDown {'scancode': 1, 'key': 27, 'unicode': '\x1b', 'mod': 0})>
<Event(3-KeyUp {'scancode': 1, 'key': 27, 'mod': 0})>
...
<Event(4-MouseMotion {'buttons': (0, 0, 0), 'pos': (323, 194), 'rel': (-3, -1)})>
<Event(4-MouseMotion {'buttons': (0, 0, 0), 'pos': (322, 193), 'rel': (-1, -1)})>
<Event(4-MouseMotion {'buttons': (0, 0, 0), 'pos': (321, 192), 'rel': (-1, -1)})>
<Event(4-MouseMotion {'buttons': (0, 0, 0), 'pos': (319, 192), 'rel': (-2, 0)})>
<Event(5-MouseButtonDown {'button': 1, 'pos': (319, 192)})>
<Event(6-MouseButtonUp {'button': 1, 'pos': (319, 192)})>
<Event(4-MouseMotion {'buttons': (0, 0, 0), 'pos': (319, 191), 'rel': (0, -1)})>
<Event(5-MouseButtonDown {'button': 2, 'pos': (319, 191)})>
<Event(5-MouseButtonDown {'button': 5, 'pos': (319, 191)})>
<Event(6-MouseButtonUp {'button': 5, 'pos': (319, 191)})>
<Event(6-MouseButtonUp {'button': 2, 'pos': (319, 191)})>
<Event(5-MouseButtonDown {'button': 3, 'pos': (319, 191)})>
<Event(6-MouseButtonUp {'button': 3, 'pos': (319, 191)})>
...
<Event(1-ActiveEvent {'state': 1, 'gain': 0})>
<Event(12-Quit {})>
```

So let us now make these changes to the code near the top of our game loop:

```
1  while True:
2
3      # Look for an event from keyboard, mouse, etc.
4      ev = pygame.event.poll()
5      if event.type == pygame.QUIT:
6          break;
7      if event.type == pygame.KEYDOWN:
8          key = ev.dict["key"]
9          if key == 27:                                # On Escape key ...
10             break                                    # leave the game loop.
11         if key == ord("r"):
12             colors[0] = (255, 0, 0)                  # Change to red + black.
13         elif key == ord("g"):
14             colors[0] = (0, 255, 0)                  # Change to green + black.
15         elif key == ord("b"):
16             colors[0] = (0, 0, 255)                  # Change to blue + black.
17
18     if event.type == pygame.MOUSEBUTTONDOWN: # Mouse gone down?
19         posn_of_click = event.dict["pos"]           # Get the coordinates.
20         print(posn_of_click)                        # Just print them.
```

Lines 7-16 show typical processing for a KEYDOWN event — if a key has gone down, we test which key it is, and take some action. With this in place, we have another way to quit our queens program — by hitting the escape key. Also, we can use keys to change the color of the board that is drawn.

Finally, at line 20, we respond (pretty lamely) to the mouse button going down.

As a final exercise in this section, we'll write a better response handler to mouse clicks. What we will do is figure out if the user has clicked the mouse on one of our sprites. If there is a sprite under the mouse when the click occurs, we'll send the click to the sprite and let it respond in some sensible way.

We'll begin with some code that finds out which sprite is under the clicked position, perhaps none! We add a method to the class, `contains_point`, which returns True if the point is within the rectangle of the sprite:

```
1  def contains_point(self, point):
2      """ Return True if my sprite rectangle contains point pt """
3      (my_x, my_y) = self.position
4      my_width = self.image.get_width()
5      my_height = self.image.get_height()
6      (x, y) = point
7      return ( x >= my_x and x < my_x + my_width and
8              y >= my_y and y < my_y + my_height)
```

Now in the game loop, once we've seen the mouse event, we determine which queen, if any, should be told to respond to the event:

```
1  if ev.type == pygame.MOUSEBUTTONDOWN:
2      posn_of_click = event.dict["pos"]
3      for sprite in all_sprites:
4          if sprite.contains_point(posn_of_click):
5              sprite.handle_click()
6              break
```

And the final thing is to write a new method called `handle_click` in the `QueenSprite` class. When a sprite is clicked, we'll just add some velocity in the up direction, i.e. kick it back into the air.

```
1  def handle_click(self):
2      self.y_velocity += -0.3  # Kick it up
```

With these changes we have a playable game! See if you can keep all the balls on the move, not allowing any one to settle!

## 14.6 A wave of animation

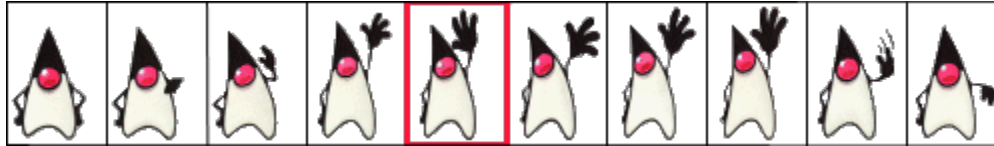
Many games have sprites that are animated: they crouch, jump and shoot. How do they do that?

Consider this sequence of 10 images: if we display them in quick succession, Duke will wave at us. (Duke is a friendly visitor from the kingdom of Javaland.)



A compound image containing smaller *patches* which are intended for animation is called a **sprite sheet**. Download this sprite sheet by right-clicking in your browser and saving it in your working directory with the name `duke_spritesheet.png`.

The sprite sheet has been quite carefully prepared: each of the 10 patches are spaced exactly 50 pixels apart. So, assuming we want to draw patch number 4 (numbering from 0), we want to draw only the rectangle that starts at x position 200, and is 50 pixels wide, within the sprite sheet. Here we’ve shown the patches and highlighted the patch we want to draw.



The `blit` method we’ve been using — for copying pixels from one surface to another — can copy a sub-rectangle of the source surface. So the grand idea here is that each time we draw Duke, we won’t blit the whole sprite sheet. Instead we’ll provide an extra rectangle argument that determines which portion of the sprite sheet will be blitted.

We’re going to add new code in this section to our existing N queens drawing game. What we want is to put some instances of Duke on the chessboard somewhere. If the user clicks on one of them, we’ll get him to respond by waving back, for one cycle of his animation.

But before we do that, we need another change. Up until now, our game loop has been running at really fast frame rates that are unpredictable. So we’ve chosen some *magic numbers* for gravity and for bouncing and kicking the ball on the basis of trial-and-error. If we’re going to start animating more sprites, we need to tame our game loop to operate at a fixed, known frame rate. This will allow us to plan our animation better.

PyGame gives us the tools to do this in just two lines of code. In the setup section of the game, we instantiate a new `Clock` object:

```
1 my_clock = pygame.time.Clock()
```

and right at the bottom of the game loop, we call a method on this object that limits the frame rate to whatever we specify. So let’s plan our game and animation for 60 frames per second, by adding this line at the bottom of our game loop:

```
1 my_clock.tick(60) # Waste time so that frame rate becomes 60 fps
```

You’ll find that you have to go back and adjust the numbers for gravity and kicking the ball now, to match this much slower frame rate. When we plan an animation so that it only works sensibly at a fixed frame rate, we say that we’ve *baked* the animation. In this case we’re baking our animations for 60 frames per second.

To fit into the existing framework that we already have for our queens board, we want to create a `DukeSprite` class that has all the same methods as the `QueenSprite` class. Then we can add one or more Duke instances onto our list of `all_sprites`, and our existing game loop will then call methods of the Duke instance. Let us start with skeleton scaffolding for the new class:

```
1 class DukeSprite:
2
3     def __init__(self, img, target_position):
4         self.image = img
5         self.position = target_position
6
7     def update(self):
8         return
9
10    def draw(self, target_surface):
11        return
12
13    def handle_click(self):
14        return
15
```

(continues on next page)



(continued from previous page)

```
16     def contains_point(self, pt):
17         # Use code from QueenSprite here
18         return
```

The only changes we'll need to the existing game are all in the setup section. We load up the new sprite sheet and instantiate a couple of instances of Duke, at the positions we want on the chessboard. So before entering the game loop, we add this code:

```
1  # Load the sprite sheet
2  duke_sprite_sheet = pygame.image.load("duke_spritesheet.png")
3
4  # Instantiate two duke instances, put them on the chessboard
5  duke1 = DukeSprite(duke_sprite_sheet, (square_size*2, 0))
6  duke2 = DukeSprite(duke_sprite_sheet, (square_size*5, sq_sz))
7
8  # Add them to the list of sprites which our game loop manages
9  all_sprites.append(duke1)
10 all_sprites.append(duke2)
```

Now the game loop will test if each instance has been clicked, will call the click handler for that instance. It will also call update and draw for all sprites. All the remaining changes we need to make will be made in the methods of the DukeSprite class.

Let's begin with drawing one of the patches. We'll introduce a new attribute `curr_patch_num` into the class. It holds a value between 0 and 9, and determines which patch to draw. So the job of the draw method is to compute the sub-rectangle of the patch to be drawn, and to blit only that portion of the spritesheet:

```
1  def draw(self, target_surface):
2      patch_rect = (self.curr_patch_num * 50, 0,
3                  50, self.image.get_height())
4      target_surface.blit(self.image, self.posn, patch_rect)
```

Now on to getting the animation to work. We need to arrange logic in update so that if we're busy animating, we change the `curr_patch_num` every so often, and we also decide when to bring Duke back to his rest position, and stop the animation. An important issue is that the game loop frame rate — in our case 60 fps — is not the same as the *animation rate* — the rate at which we want to change Duke's animation patches. So we'll plan Duke wave's animation cycle for a duration of 1 second. In other words, we want to play out Duke's 10 animation patches over 60 calls to update. (This is how the baking of the animation takes place!) So we'll keep another animation frame counter in the class, which will be zero when we're not animating, and each call to update will increment the counter up to 59, and then back to 0. We can then divide that animation counter by 6, to set the `curr_patch_num` variable to select the patch we want to show.

```
1  def update(self):
2      if self.anim_frame_count > 0:
3          self.anim_frame_count = (self.anim_frame_count + 1) % 60
4          self.curr_patch_num = self.anim_frame_count // 6
```

Notice that if `anim_frame_count` is zero, i.e. Duke is at rest, nothing happens here. But if we start the counter running, it will count up to 59 before settling back to zero. Notice also, that because `anim_frame_count` can only be a value between 0 and 59, the `curr_patch_num` will always stay between 0 and 9. Just what we require!

Now how do we trigger the animation, and start it running? On the mouse click.

```
1  def handle_click(self):
2      if self.anim_frame_count == 0:
3          self.anim_frame_count = 5
```

Two things of interest here. We only start the animation if Duke is at rest. Clicks on Duke while he is already waving get ignored. And when we do start the animation, we set the counter to 5 — this means that on the very next call to `update` the counter becomes 6, and the image changes. If we had set the counter to 1, we would have needed to wait for 5 more calls to `update` before anything happened — a slight lag, but enough to make things feel sluggish.

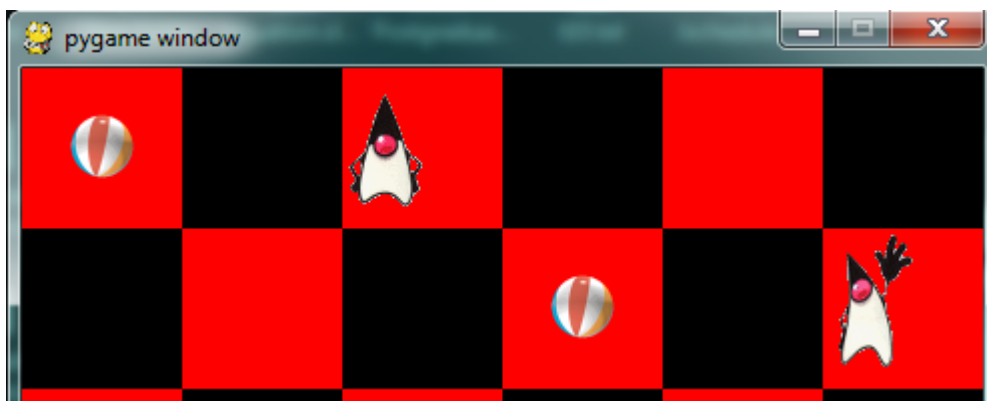
The final touch-up is to initialize our two new attributes when we instantiate the class. Here is the code for the whole class now:

```

1  class DukeSprite:
2
3      def __init__(self, img, target_posn):
4          self.image = img
5          self.position = target_posn
6          self.anim_frame_count = 0
7          self.curr_patch_num = 0
8
9      def update(self):
10         if self.anim_frame_count > 0:
11             self.anim_frame_count = (self.anim_frame_count + 1) % 60
12             self.curr_patch_num = self.anim_frame_count // 6
13
14         def draw(self, target_surface):
15             patch_rect = (self.curr_patch_num * 50, 0,
16                           50, self.image.get_height())
17             target_surface.blit(self.image, self.posn, patch_rect)
18
19         def contains_point(self, pt):
20             """ Return True if my sprite rectangle contains pt """
21             (my_x, my_y) = self.posn
22             my_width = self.image.get_width()
23             my_height = self.image.get_height()
24             (x, y) = pt
25             return ( x >= my_x and x < my_x + my_width and
26                     y >= my_y and y < my_y + my_height)
27
28         def handle_click(self):
29             if self.anim_frame_count == 0:
30                 self.anim_frame_count = 5

```

Now we have two extra Duke instances on our chessboard, and clicking on either causes that instance to wave.



## 14.7 Aliens - a case study

Find the example games with the PyGame package, (On a windows system, something like C:\Python3\Lib\site-packages\pygame\examples) and play the Aliens game. Then read the code, in an editor or Python environment that shows line numbers.

It does a number of much more advanced things that we do, and relies on the PyGame framework for more of its logic. Here are some of the points to notice:

- The frame rate is deliberately constrained near the bottom of the game loop at line 311. If we change that number we can make the game very slow or unplayably fast!
- There are different kinds of sprites: Explosions, Shots, Bombs, Aliens and a Player. Some of these have more than one image — by swapping the images, we get animation of the sprites, i.e. the Alien spacecraft lights change, and this is done at line 112.
- Different kinds of objects are referenced in different groups of sprites, and PyGame helps maintain these. This lets the program check for collisions between, say, the list of shots fired by the player, and the list of spaceships that are attacking. PyGame does a lot of the hard work for us.
- Unlike our game, objects in the Aliens game have a limited lifetime, and have to get killed. For example, if we shoot, a Shot object is created — if it reaches the top of the screen without exploding against anything, it has to be removed from the game. Lines 141-142 do this. Similarly, when a falling bomb gets close to the ground (line 156), it instantiates a new Explosion sprite, and the bomb kills itself.
- There are random timings that add to the fun — when to spawn the next Alien, when an Alien drops the next bomb, etc.
- The game plays sounds too: a less-than-relaxing loop sound, plus sounds for the shots and explosions.

## 14.8 Reflections

Object oriented programming is a good organizational tool for software. In the examples in this chapter, we've started to use (and hopefully appreciate) these benefits. Here we had N queens each with its own state, falling to its own floor level, bouncing, getting kicked, etc. We might have managed without the organizational power of objects — perhaps we could have kept lists of velocities for each queen, and lists of target positions, and so on — our code would likely have been much more complicated, ugly, and a lot poorer!

## 14.9 Glossary

**animation rate** The rate at which we play back successive patches to create the illusion of movement. In the sample we considered in this chapter, we played Duke's 10 patches over the duration of one second. Not the same as the frame rate.

**baked animation** An animation that is designed to look good at a predetermined fixed frame rate. This reduces the amount of computation that needs to be done when the game is running. High-end commercial games usually bake their animations.

**blit** A verb used in computer graphics, meaning to make a fast copy of an image or pixels from a sub-rectangle of one image or surface to another surface or image.

**frame rate** The rate at which the game loop executes and updates the display.

**game loop** A loop that drives the logic of a game. It will usually poll for events, then update each of the objects in the game, then get everything drawn, and then put the newly drawn frame on display.

**pixel** A single picture element, or dot, from which images are made.

**poll** To ask whether something like a keypress or mouse movement has happened. Game loops usually poll to discover what events have occurred. This is different from event-driven programs like the ones seen in the chapter titled “Events”. In those cases, the button click or keypress event triggers the call of a handler function in your program, but this happens behind your back.

**sprite** An active agent or element in a game, with its own state, position and behaviour.

**surface** This is PyGame’s term for what the Turtle module calls a *canvas*. A surface is a rectangle of pixels used for displaying shapes and images.

## 14.10 Exercises

1. Have fun with Python, and with PyGame.
2. We deliberately left a bug in the code for animating Duke. If you click on one of the chessboard squares to the right of Duke, he waves anyway. Why? Find a one-line fix for the bug.
3. Use your preferred search engine to search their image library for “sprite sheet playing cards”. Create a list [0..51] to represent an encoding of the 52 cards in a deck. Shuffle the cards, slice off the top five as your hand in a poker deal. Display the hand you have been dealt.
4. So the Aliens game is in outer space, without gravity. Shots fly away forever, and bombs don’t speed up when they fall. Add some gravity to the game. Decide if you’re going to allow your own shots to fall back on your head and kill you.
5. Those pesky Aliens seem to pass right through each other! Change the game so that they collide, and destroy each other in a mighty explosion.

`while` loop, within another `if` statement, and so on.

Having seen this pattern, and having learned about lists and objects, we should not be surprised to learn that we can create lists of objects. We can also create objects that contain lists (as attributes); we can create lists that contain lists; we can create objects that contain objects; and so on.

In this chapter and the next, we will look at some examples of these combinations, using `Card` objects as an example.

### D.4.2 Card objects

If you are not familiar with common playing cards, now would be a good time to get a deck, or else this chapter might not make much sense. There are fifty-two cards in a deck, each of which belongs to one of four suits and one of thirteen ranks. The suits are Spades, Hearts, Diamonds, and Clubs (in descending order in bridge). The ranks are Ace, 2, 3, 4, 5, 6, 7, 8, 9, 10, Jack, Queen, and King. Depending on the game that we are playing, the rank of Ace may be higher than King or lower than 2. The rank is sometimes called the face-value of the card.

If we want to define a new object to represent a playing card, it is obvious what the attributes should be: `rank` and `suit`. It is not as obvious what type the attributes should be. One possibility is to use strings containing words like "Spade" for suits and "Queen" for ranks. One problem with this implementation is that it would not be easy to compare cards to see which had a higher rank or suit.

An alternative is to use integers to **encode** the ranks and suits. By encode, we do not mean what some people think, which is to encrypt or translate into a secret code. What a computer scientist means by encode is to define a mapping between a sequence of numbers and the items I want to represent. For example:

```
Spades    -->  3
Hearts    -->  2
Diamonds  -->  1
Clubs     -->  0
```

An obvious feature of this mapping is that the suits map to integers in order, so we can compare suits by comparing integers. The mapping for ranks is fairly obvious; each of the numerical ranks maps to the corresponding integer, and for face cards:

```
Jack      --> 11
Queen     --> 12
King      --> 13
```

The reason we are using mathematical notation for these mappings is that they are not part of the Python program. They are part of the program design, but they never appear explicitly in the code. The class definition for the `Card` type looks like this:

```
1 class Card:
2     def __init__(self, suit=0, rank=0):
3         self.suit = suit
4         self.rank = rank
```

As usual, we provide an initialization method that takes an optional parameter for each attribute.

To create some objects, representing say the 3 of Clubs and the Jack of Diamonds, use these commands:

```
1 three_of_clubs = Card(0, 3)
2 card1 = Card(1, 11)
```

In the first case above, for example, the first argument, 0, represents the suit Clubs.

---

Save this code for later use ...

In the next chapter we assume that we have save the `Cards` class, and the upcoming `Deck` class in a file called `Cards.py`.

---

### D.4.3 Class attributes and the `__str__` method

In order to print `Card` objects in a way that people can easily read, we want to map the integer codes onto words. A natural way to do that is with lists of strings. We assign these lists to **class attributes** at the top of the class definition:

```
1 class Card:
2     suits = ["Clubs", "Diamonds", "Hearts", "Spades"]
3     ranks = ["narf", "Ace", "2", "3", "4", "5", "6", "7",
4             "8", "9", "10", "Jack", "Queen", "King"]
5
6     def __init__(self, suit=0, rank=0):
7         self.suit = suit
8         self.rank = rank
9
10    def __str__(self):
11        return (self.ranks[self.rank] + " of " + self.suits[self.suit])
```

A class attribute is defined outside of any method, and it can be accessed from any of the methods in the class.

Inside `__str__`, we can use `suits` and `ranks` to map the numerical values of `suit` and `rank` to strings. For example, the expression `self.suits[self.suit]` means use the attribute `suit` from the object `self` as an index into the class attribute named `suits`, and select the appropriate string.

The reason for the "narf" in the first element in `ranks` is to act as a place keeper for the zero-eth element of the list, which will never be used. The only valid ranks are 1 to 13. This wasted item is not entirely necessary. We could have started at 0, as usual, but it is less confusing to encode the rank 2 as integer 2, 3 as 3, and so on.

With the methods we have so far, we can create and print cards:

```
>>> card1 = Card(1, 11)
>>> print(card1)
Jack of Diamonds
```

Class attributes like `suits` are shared by all `Card` objects. The advantage of this is that we can use any `Card` object to access the class attributes:

```
>>> card2 = Card(1, 3)
>>> print(card2)
3 of Diamonds
>>> print(card2.suits[1])
Diamonds
```

Because every `Card` instance references the same class attribute, we have an aliasing situation. The disadvantage is that if we modify a class attribute, it affects every instance of the class. For example, if we decide that Jack of Diamonds should really be called Jack of Swirly Whales, we could do this:

```
>>> card1.suits[1] = "Swirly Whales"
>>> print(card1)
Jack of Swirly Whales
```

The problem is that *all* of the Diamonds just became Swirly Whales:

```
>>> print(card2)
3 of Swirly Whales
```

It is usually not a good idea to modify class attributes.

### D.4.4 Comparing cards

For primitive types, there are six relational operators (<, >, ==, etc.) that compare values and determine when one is greater than, less than, or equal to another. If we want our own types to be comparable using the syntax of these relational operators, we need to define six corresponding special methods in our class.

We'd like to start with a single method named `cmp` that houses the logic of ordering. By convention, a comparison method takes two parameters, `self` and `other`, and returns 1 if the first object is greater, -1 if the second object is greater, and 0 if they are equal to each other.

Some types are completely ordered, which means that we can compare any two elements and tell which is bigger. For example, the integers and the floating-point numbers are completely ordered. Some types are unordered, which means that there is no meaningful way to say that one element is bigger than another. For example, the fruits are unordered, which is why we cannot compare apples and oranges, and we cannot meaningfully order a collection of images, or a collection of cellphones.

Playing cards are partially ordered, which means that sometimes we can compare cards and sometimes not. For example, we know that the 3 of Clubs is higher than the 2 of Clubs, and the 3 of Diamonds is higher than the 3 of Clubs. But which is better, the 3 of Clubs or the 2 of Diamonds? One has a higher rank, but the other has a higher suit.

In order to make cards comparable, we have to decide which is more important, rank or suit. To be honest, the choice is arbitrary. For the sake of choosing, we will say that suit is more important, because a new deck of cards comes sorted with all the Clubs together, followed by all the Diamonds, and so on.

With that decided, we can write `cmp`:

```
1 def cmp(self, other):
2     # Check the suits
3     if self.suit > other.suit: return 1
4     if self.suit < other.suit: return -1
5     # Suits are the same... check ranks
6     if self.rank > other.rank: return 1
7     if self.rank < other.rank: return -1
8     # Ranks are the same... it's a tie
9     return 0
```

In this ordering, Aces appear lower than Deuces (2s).

Now, we can define the six special methods that do the overloading of each of the relational operators for us:

```
1 def __eq__(self, other):
2     return self.cmp(other) == 0
3
4 def __le__(self, other):
5     return self.cmp(other) <= 0
6
7 def __ge__(self, other):
8     return self.cmp(other) >= 0
9
10 def __gt__(self, other):
11     return self.cmp(other) > 0
12
```

(continues on next page)

(continued from previous page)

```
13 def __lt__(self, other):
14     return self.cmp(other) < 0
15
16 def __ne__(self, other):
17     return self.cmp(other) != 0
```

With this machinery in place, the relational operators now work as we'd like them to:

```
>>> card1 = Card(1, 11)
>>> card2 = Card(1, 3)
>>> card3 = Card(1, 11)
>>> card1 < card2
False
>>> card1 == card3
True
```

## D.4.5 Decks

Now that we have objects to represent Cards, the next logical step is to define a class to represent a Deck. Of course, a deck is made up of cards, so each Deck object will contain a list of cards as an attribute. Many card games will need at least two different decks — a red deck and a blue deck.

The following is a class definition for the Deck class. The initialization method creates the attribute `cards` and generates the standard pack of fifty-two cards:

```
1 class Deck:
2     def __init__(self):
3         self.cards = []
4         for suit in range(4):
5             for rank in range(1, 14):
6                 self.cards.append(Card(suit, rank))
```

The easiest way to populate the deck is with a nested loop. The outer loop enumerates the suits from 0 to 3. The inner loop enumerates the ranks from 1 to 13. Since the outer loop iterates four times, and the inner loop iterates thirteen times, the total number of times the body is executed is fifty-two (thirteen times four). Each iteration creates a new instance of `Card` with the current suit and rank, and appends that card to the `cards` list.

With this in place, we can instantiate some decks:

```
1 red_deck = Deck()
2 blue_deck = Deck()
```

## D.4.6 Printing the deck

As usual, when we define a new type we want a method that prints the contents of an instance. To print a Deck, we traverse the list and print each Card:

```
1 class Deck:
2     ...
3     def print_deck(self):
4         for card in self.cards:
5             print(card)
```



Here, and from now on, the ellipsis ( . . . ) indicates that we have omitted the other methods in the class.

As an alternative to `print_deck`, we could write a `__str__` method for the `Deck` class. The advantage of `__str__` is that it is more flexible. Rather than just printing the contents of the object, it generates a string representation that other parts of the program can manipulate before printing, or store for later use.

Here is a version of `__str__` that returns a string representation of a `Deck`. To add a bit of pizzazz, it arranges the cards in a cascade where each card is indented one space more than the previous card:

```
1 class Deck:
2     ...
3     def __str__(self):
4         s = ""
5         for i in range(len(self.cards)):
6             s = s + " " * i + str(self.cards[i]) + "\n"
7         return s
```

This example demonstrates several features. First, instead of traversing `self.cards` and assigning each card to a variable, we are using `i` as a loop variable and an index into the list of cards.

Second, we are using the string multiplication operator to indent each card by one more space than the last. The expression `" " * i` yields a number of spaces equal to the current value of `i`.

Third, instead of using the `print` command to print the cards, we use the `str` function. Passing an object as an argument to `str` is equivalent to invoking the `__str__` method on the object.

Finally, we are using the variable `s` as an **accumulator**. Initially, `s` is the empty string. Each time through the loop, a new string is generated and concatenated with the old value of `s` to get the new value. When the loop ends, `s` contains the complete string representation of the `Deck`, which looks like this:

```
>>> red_deck = Deck()
>>> print(red_deck)
Ace of Clubs
 2 of Clubs
 3 of Clubs
 4 of Clubs
 5 of Clubs
 6 of Clubs
 7 of Clubs
 8 of Clubs
 9 of Clubs
10 of Clubs
 Jack of Clubs
  Queen of Clubs
   King of Clubs
    Ace of Diamonds
     2 of Diamonds
     ...
```

And so on. Even though the result appears on 52 lines, it is one long string that contains newlines.

### D.4.7 Shuffling the deck

If a deck is perfectly shuffled, then any card is equally likely to appear anywhere in the deck, and any location in the deck is equally likely to contain any card.

To shuffle the deck, we will use the `randrange` function from the `random` module. With two integer arguments, `a` and `b`, `randrange` chooses a random integer in the range `a <= x < b`. Since the upper bound is strictly less than `b`, we can use the length of a list as the second parameter, and we are guaranteed to get a legal index. For example, if

`rng` has already been instantiated as a random number source, this expression chooses the index of a random card in a deck:

```
1 rng.randrange(0, len(self.cards))
```

An easy way to shuffle the deck is by traversing the cards and swapping each card with a randomly chosen one. It is possible that the card will be swapped with itself, but that is fine. In fact, if we precluded that possibility, the order of the cards would be less than entirely random:

```
1 class Deck:
2     ...
3     def shuffle(self):
4         import random
5         rng = random.Random()           # Create a random generator
6         num_cards = len(self.cards)
7         for i in range(num_cards):
8             j = rng.randrange(i, num_cards)
9             (self.cards[i], self.cards[j]) = (self.cards[j], self.cards[i])
```

Rather than assume that there are fifty-two cards in the deck, we get the actual length of the list and store it in `num_cards`.

For each card in the deck, we choose a random card from among the cards that haven't been shuffled yet. Then we swap the current card (`i`) with the selected card (`j`). To swap the cards we use a tuple assignment:

```
1 (self.cards[i], self.cards[j]) = (self.cards[j], self.cards[i])
```

While this is a good shuffling method, a random number generator object also has a `shuffle` method that can shuffle elements in a list, in place. So we could rewrite this function to use the one provided for us:

```
1 class Deck:
2     ...
3     def shuffle(self):
4         import random
5         rng = random.Random()           # Create a random generator
6         rng.shuffle(self.cards)         # uUse its shuffle method
```

## D.4.8 Removing and dealing cards

Another method that would be useful for the `Deck` class is `remove`, which takes a card as a parameter, removes it, and returns `True` if the card was in the deck and `False` otherwise:

```
1 class Deck:
2     ...
3     def remove(self, card):
4         if card in self.cards:
5             self.cards.remove(card)
6             return True
7         else:
8             return False
```

The `in` operator returns `True` if the first operand is in the second. If the first operand is an object, Python uses the object's `__eq__` method to determine equality with items in the list. Since the `__eq__` we provided in the `Card` class checks for deep equality, the `remove` method checks for deep equality.

To deal cards, we want to remove and return the top card. The list method `pop` provides a convenient way to do that:

```
1 class Deck:
2     ...
3     def pop(self):
4         return self.cards.pop()
```

Actually, `pop` removes the *last* card in the list, so we are in effect dealing from the bottom of the deck.

One more operation that we are likely to want is the Boolean function `is_empty`, which returns `True` if the deck contains no cards:

```
1 class Deck:
2     ...
3     def is_empty(self):
4         return self.cards == []
```

## D.4.9 Glossary

**encode** To represent one type of value using another type of value by constructing a mapping between them.

**class attribute** A variable that is defined inside a class definition but outside any method. Class attributes are accessible from any method in the class and are shared by all instances of the class.

**accumulator** A variable used in a loop to accumulate a series of values, such as by concatenating them onto a string or adding them to a running sum.

## D.4.10 Exercises

1. Modify `cmp` so that Aces are ranked higher than Kings.

## D.5 Inheritance

### D.5.1 Inheritance

The language feature most often associated with object-oriented programming is **inheritance**. Inheritance is the ability to define a new class that is a modified version of an existing class.

The primary advantage of this feature is that you can add new methods to a class without modifying the existing class. It is called inheritance because the new class inherits all of the methods of the existing class. Extending this metaphor, the existing class is sometimes called the **parent** class. The new class may be called the **child** class or sometimes subclass.

Inheritance is a powerful feature. Some programs that would be complicated without inheritance can be written concisely and simply with it. Also, inheritance can facilitate code reuse, since you can customize the behavior of parent classes without having to modify them. In some cases, the inheritance structure reflects the natural structure of the problem, which makes the program easier to understand.

On the other hand, inheritance can make programs difficult to read. When a method is invoked, it is sometimes not clear where to find its definition. The relevant code may be scattered among several modules. Also, many of the things that can be done using inheritance can be done as elegantly (or more so) without it. If the natural structure of the problem does not lend itself to inheritance, this style of programming can do more harm than good.

In this chapter we will demonstrate the use of inheritance as part of a program that plays the card game Old Maid. One of our goals is to write code that could be reused to implement other card games.

### D.5.2 A hand of cards

For almost any card game, we need to represent a hand of cards. A hand is similar to a deck, of course. Both are made up of a set of cards, and both require operations like adding and removing cards. Also, we might like the ability to shuffle both decks and hands.

A hand is also different from a deck. Depending on the game being played, we might want to perform some operations on hands that don't make sense for a deck. For example, in poker we might classify a hand (straight, flush, etc.) or compare it with another hand. In bridge, we might want to compute a score for a hand in order to make a bid.

This situation suggests the use of inheritance. If `Hand` is a subclass of `Deck`, it will have all the methods of `Deck`, and new methods can be added.

We add the code in this chapter to our `Cards.py` file from the previous chapter. In the class definition, the name of the parent class appears in parentheses:

```
1 class Hand(Deck):
2     pass
```

This statement indicates that the new `Hand` class inherits from the existing `Deck` class.

The `Hand` constructor initializes the attributes for the hand, which are `name` and `cards`. The string `name` identifies this hand, probably by the name of the player that holds it. The name is an optional parameter with the empty string as a default value. `cards` is the list of cards in the hand, initialized to the empty list:

```
1 class Hand(Deck):
2     def __init__(self, name=""):
3         self.cards = []
4         self.name = name
```

For just about any card game, it is necessary to add and remove cards from the deck. Removing cards is already taken care of, since `Hand` inherits `remove` from `Deck`. But we have to write `add`:

```
1 class Hand(Deck):
2     ...
3     def add(self, card):
4         self.cards.append(card)
```

Again, the ellipsis indicates that we have omitted other methods. The list `append` method adds the new card to the end of the list of cards.

### D.5.3 Dealing cards

Now that we have a `Hand` class, we want to deal cards from the `Deck` into hands. It is not immediately obvious whether this method should go in the `Hand` class or in the `Deck` class, but since it operates on a single deck and (possibly) several hands, it is more natural to put it in `Deck`.

`deal` should be fairly general, since different games will have different requirements. We may want to deal out the entire deck at once or add one card to each hand.

`deal` takes two parameters, a list (or tuple) of hands and the total number of cards to deal. If there are not enough cards in the deck, the method deals out all of the cards and stops:

```
1 class Deck:
2     ...
3     def deal(self, hands, num_cards=999):
4         num_hands = len(hands)
5         for i in range(num_cards):
6             if self.is_empty():
7                 break                # Break if out of cards
8             card = self.pop()        # Take the top card
9             hand = hands[i % num_hands] # Whose turn is next?
10            hand.add(card)           # Add the card to the hand
```

The second parameter, `num_cards`, is optional; the default is a large number, which effectively means that all of the cards in the deck will get dealt.

The loop variable `i` goes from 0 to `num_cards-1`. Each time through the loop, a card is removed from the deck using the list method `pop`, which removes and returns the last item in the list.

The modulus operator (`%`) allows us to deal cards in a round robin (one card at a time to each hand). When `i` is equal to the number of hands in the list, the expression `i % num_hands` wraps around to the beginning of the list (index 0).

## D.5.4 Printing a Hand

To print the contents of a hand, we can take advantage of the `__str__` method inherited from `Deck`. For example:

```
>>> deck = Deck()
>>> deck.shuffle()
>>> hand = Hand("frank")
>>> deck.deal([hand], 5)
>>> print(hand)
Hand frank contains
2 of Spades
3 of Spades
4 of Spades
Ace of Hearts
9 of Clubs
```

It's not a great hand, but it has the makings of a straight flush.

Although it is convenient to inherit the existing methods, there is additional information in a `Hand` object we might want to include when we print one. To do that, we can provide a `__str__` method in the `Hand` class that overrides the one in the `Deck` class:

```
1 class Hand(Deck)
2     ...
3     def __str__(self):
4         s = "Hand " + self.name
5         if self.is_empty():
6             s += " is empty\n"
7         else:
8             s += " contains\n"
9         return s + Deck.__str__(self)
```

Initially, `s` is a string that identifies the hand. If the hand is empty, the program appends the words `is empty` and returns `s`.

Otherwise, the program appends the word `contains` and the string representation of the `Deck`, computed by invoking the `__str__` method in the `Deck` class on `self`.

It may seem odd to send `self`, which refers to the current `Hand`, to a `Deck` method, until you remember that a `Hand` is a kind of `Deck`. `Hand` objects can do everything `Deck` objects can, so it is legal to send a `Hand` to a `Deck` method.

In general, it is always legal to use an instance of a subclass in place of an instance of a parent class.

### D.5.5 The `CardGame` class

The `CardGame` class takes care of some basic chores common to all games, such as creating the deck and shuffling it:

```
1 class CardGame:
2     def __init__(self):
3         self.deck = Deck()
4         self.deck.shuffle()
```

This is the first case we have seen where the initialization method performs a significant computation, beyond initializing attributes.

To implement specific games, we can inherit from `CardGame` and add features for the new game. As an example, we'll write a simulation of Old Maid.

The object of Old Maid is to get rid of cards in your hand. You do this by matching cards by rank and color. For example, the 4 of Clubs matches the 4 of Spades since both suits are black. The Jack of Hearts matches the Jack of Diamonds since both are red.

To begin the game, the Queen of Clubs is removed from the deck so that the Queen of Spades has no match. The fifty-one remaining cards are dealt to the players in a round robin. After the deal, all players match and discard as many cards as possible.

When no more matches can be made, play begins. In turn, each player picks a card (without looking) from the closest neighbor to the left who still has cards. If the chosen card matches a card in the player's hand, the pair is removed. Otherwise, the card is added to the player's hand. Eventually all possible matches are made, leaving only the Queen of Spades in the loser's hand.

In our computer simulation of the game, the computer plays all hands. Unfortunately, some nuances of the real game are lost. In a real game, the player with the Old Maid goes to some effort to get their neighbor to pick that card, by displaying it a little more prominently, or perhaps failing to display it more prominently, or even failing to fail to display that card more prominently. The computer simply picks a neighbor's card at random.

### D.5.6 `OldMaidHand` class

A hand for playing Old Maid requires some abilities beyond the general abilities of a `Hand`. We will define a new class, `OldMaidHand`, that inherits from `Hand` and provides an additional method called `remove_matches`:

```
1 class OldMaidHand(Hand):
2     def remove_matches(self):
3         count = 0
4         original_cards = self.cards[:]
5         for card in original_cards:
6             match = Card(3 - card.suit, card.rank)
7             if match in self.cards:
8                 self.cards.remove(card)
9                 self.cards.remove(match)
10                print("Hand {0}: {1} matches {2}"
11                      .format(self.name, card, match))
12                count += 1
13        return count
```

We start by making a copy of the list of cards, so that we can traverse the copy while removing cards from the original. Since `self.cards` is modified in the loop, we don't want to use it to control the traversal. Python can get quite confused if it is traversing a list that is changing!

For each card in the hand, we figure out what the matching card is and go looking for it. The match card has the same rank and the other suit of the same color. The expression `3 - card.suit` turns a Club (suit 0) into a Spade (suit 3) and a Diamond (suit 1) into a Heart (suit 2). You should satisfy yourself that the opposite operations also work. If the match card is also in the hand, both cards are removed.

The following example demonstrates how to use `remove_matches`:

```
>>> game = CardGame()
>>> hand = OldMaidHand("frank")
>>> game.deck.deal([hand], 13)
>>> print(hand)
Hand frank contains
Ace of Spades
2 of Diamonds
7 of Spades
8 of Clubs
6 of Hearts
8 of Spades
7 of Clubs
Queen of Clubs
7 of Diamonds
5 of Clubs
Jack of Diamonds
10 of Diamonds
10 of Hearts
>>> hand.remove_matches()
Hand frank: 7 of Spades matches 7 of Clubs
Hand frank: 8 of Spades matches 8 of Clubs
Hand frank: 10 of Diamonds matches 10 of Hearts
>>> print(hand)
Hand frank contains
Ace of Spades
2 of Diamonds
6 of Hearts
Queen of Clubs
7 of Diamonds
5 of Clubs
Jack of Diamonds
```

Notice that there is no `__init__` method for the `OldMaidHand` class. We inherit it from `Hand`.

### D.5.7 OldMaidGame class

Now we can turn our attention to the game itself. `OldMaidGame` is a subclass of `CardGame` with a new method called `play` that takes a list of players as a parameter.

Since `__init__` is inherited from `CardGame`, a new `OldMaidGame` object contains a new shuffled deck:

```
1 class OldMaidGame(CardGame):
2     def play(self, names):
3         # Remove Queen of Clubs
4         self.deck.remove(Card(0,12))
5
```

(continues on next page)

(continued from previous page)

```

6         # Make a hand for each player
7         self.hands = []
8         for name in names:
9             self.hands.append(OldMaidHand(name))
10
11        # Deal the cards
12        self.deck.deal(self.hands)
13        print("----- Cards have been dealt")
14        self.print_hands()
15
16        # Remove initial matches
17        matches = self.remove_all_matches()
18        print("----- Matches discarded, play begins")
19        self.print_hands()
20
21        # Play until all 50 cards are matched
22        turn = 0
23        num_hands = len(self.hands)
24        while matches < 25:
25            matches += self.play_one_turn(turn)
26            turn = (turn + 1) % num_hands
27
28        print("----- Game is Over")
29        self.print_hands()

```

The writing of `print_hands` has been left as an exercise.

Some of the steps of the game have been separated into methods. `remove_all_matches` traverses the list of hands and invokes `remove_matches` on each:

```

1 class OldMaidGame(CardGame):
2     ...
3     def remove_all_matches(self):
4         count = 0
5         for hand in self.hands:
6             count += hand.remove_matches()
7         return count

```

`count` is an accumulator that adds up the number of matches in each hand. When we've gone through every hand, the total is returned (`count`).

When the total number of matches reaches twenty-five, fifty cards have been removed from the hands, which means that only one card is left and the game is over.

The variable `turn` keeps track of which player's turn it is. It starts at 0 and increases by one each time; when it reaches `num_hands`, the modulus operator wraps it back around to 0.

The method `play_one_turn` takes a parameter that indicates whose turn it is. The return value is the number of matches made during this turn:

```

1 class OldMaidGame(CardGame):
2     ...
3     def play_one_turn(self, i):
4         if self.hands[i].is_empty():
5             return 0
6         neighbor = self.find_neighbor(i)
7         picked_card = self.hands[neighbor].pop()
8         self.hands[i].add(picked_card)

```

(continues on next page)



(continued from previous page)

```
9         print("Hand", self.hands[i].name, "picked", picked_card)
10        count = self.hands[i].remove_matches()
11        self.hands[i].shuffle()
12        return count
```

If a player's hand is empty, that player is out of the game, so he or she does nothing and returns 0.

Otherwise, a turn consists of finding the first player on the left that has cards, taking one card from the neighbor, and checking for matches. Before returning, the cards in the hand are shuffled so that the next player's choice is random.

The method `find_neighbor` starts with the player to the immediate left and continues around the circle until it finds a player that still has cards:

```
1 class OldMaidGame(CardGame):
2     ...
3     def find_neighbor(self, i):
4         num_hands = len(self.hands)
5         for next in range(1, num_hands):
6             neighbor = (i + next) % num_hands
7             if not self.hands[neighbor].is_empty():
8                 return neighbor
```

If `find_neighbor` ever went all the way around the circle without finding cards, it would return `None` and cause an error elsewhere in the program. Fortunately, we can prove that that will never happen (as long as the end of the game is detected correctly).

We have omitted the `print_hands` method. You can write that one yourself.

The following output is from a truncated form of the game where only the top fifteen cards (tens and higher) were dealt to three players. With this small deck, play stops after seven matches instead of twenty-five.

```
>>> import cards
>>> game = cards.OldMaidGame()
>>> game.play(["Allen", "Jeff", "Chris"])
----- Cards have been dealt
Hand Allen contains
King of Hearts
Jack of Clubs
Queen of Spades
King of Spades
10 of Diamonds

Hand Jeff contains
Queen of Hearts
Jack of Spades
Jack of Hearts
King of Diamonds
Queen of Diamonds

Hand Chris contains
Jack of Diamonds
King of Clubs
10 of Spades
10 of Hearts
10 of Clubs

Hand Jeff: Queen of Hearts matches Queen of Diamonds
Hand Chris: 10 of Spades matches 10 of Clubs
```

(continues on next page)

(continued from previous page)

```
----- Matches discarded, play begins
Hand Allen contains
King of Hearts
  Jack of Clubs
    Queen of Spades
      King of Spades
        10 of Diamonds

Hand Jeff contains
Jack of Spades
  Jack of Hearts
    King of Diamonds

Hand Chris contains
Jack of Diamonds
  King of Clubs
    10 of Hearts

Hand Allen picked King of Diamonds
Hand Allen: King of Hearts matches King of Diamonds
Hand Jeff picked 10 of Hearts
Hand Chris picked Jack of Clubs
Hand Allen picked Jack of Hearts
Hand Jeff picked Jack of Diamonds
Hand Chris picked Queen of Spades
Hand Allen picked Jack of Diamonds
Hand Allen: Jack of Hearts matches Jack of Diamonds
Hand Jeff picked King of Clubs
Hand Chris picked King of Spades
Hand Allen picked 10 of Hearts
Hand Allen: 10 of Diamonds matches 10 of Hearts
Hand Jeff picked Queen of Spades
Hand Chris picked Jack of Spades
Hand Chris: Jack of Clubs matches Jack of Spades
Hand Jeff picked King of Spades
Hand Jeff: King of Clubs matches King of Spades
----- Game is Over
Hand Allen is empty

Hand Jeff contains
Queen of Spades

Hand Chris is empty
```

So Jeff loses.

## D.5.8 Glossary

**inheritance** The ability to define a new class that is a modified version of a previously defined class.

**parent class** The class from which a child class inherits.

**child class** A new class created by inheriting from an existing class; also called a subclass.

### D.5.9 Exercises

1. Add a method, `print_hands`, to the `OldMaidGame` class which traverses `self.hands` and prints each hand.
2. Define a new kind of Turtle, `TurtleGTX`, that comes with some extra features: it can jump forward a given distance, and it has an odometer that keeps track of how far the turtle has travelled since it came off the production line. (The parent class has a number of synonyms like `fd`, `forward`, `back`, `backward`, and `bk`: for this exercise, just focus on putting this functionality into the `forward` method.) Think carefully about how to count the distance if the turtle is asked to move forward by a negative amount. (We would not want to buy a second-hand turtle whose odometer reading was faked because its previous owner drove it backwards around the block too often. Try this in a car near you, and see if the car's odometer counts up or down when you reverse.)
3. After travelling some random distance, your turtle should break down with a flat tyre. After this happens, raise an exception whenever `forward` is called. Also provide a `change_tyre` method that can fix the flat.

## E.1 Catching exceptions

Whenever a runtime error occurs, it creates an **exception** object. The program stops running at this point and Python prints out the traceback, which ends with a message describing the exception that occurred.

For example, dividing by zero creates an exception:

```
>>> print(55/0)
Traceback (most recent call last):
  File "<interactive input>", line 1, in <module>
ZeroDivisionError: integer division or modulo by zero
```

So does accessing a non-existent list item:

```
>>> a = []
>>> print(a[5])
Traceback (most recent call last):
  File "<interactive input>", line 1, in <module>
IndexError: list index out of range
```

Or trying to make an item assignment on a tuple:

```
>>> tup = ("a", "b", "d", "d")
>>> tup[2] = "c"
Traceback (most recent call last):
  File "<interactive input>", line 1, in <module>
TypeError: 'tuple' object does not support item assignment
```

In each case, the error message on the last line has two parts: the type of error before the colon, and specifics about the error after the colon.

Sometimes we want to execute an operation that might cause an exception, but we don't want the program to stop. We can **handle the exception** using the `try` statement to “wrap” a region of code.

For example, we might prompt the user for the name of a file and then try to open it. If the file doesn't exist, we don't want the program to crash; we want to handle the exception:

```
1 filename = input("Enter a file name: ")
2 try:
3     f = open(filename, "r")
4 except FileNotFoundError:
5     print("There is no file named", filename)
```

The `try` statement has four separate clauses—or parts—introduced by the keywords `try`, `except`, `else`, and `finally`. All clauses but the `try` can be omitted.

The interpreter executes the block under the `try` statement, and monitors for exceptions. If one occurs, the interpreter moves to the `except` statement; it executes the `except` block if the exception raised match the exception requested in the `except` statement. If no exception occurs, the interpreter skips the block under the `except` clause. A `else` block is executed after the `try` one, if no exception occurred. A `finally` block is executed in any case. With all the statements, a `try` clause looks like:

```
1 user_input = input('Type a number:')
2 try:
3     # Try do do something that could fail.
4     user_input_as_number = float(user_input)
5 except ValueError:
6     # This will be executed if a ``ValueError`` is raised.
7     print('You did not enter a number.')
8 else:
9     # This will be executed if not exception got raised in the
10    # ``try`` statement.
11    print('The square of your number is ', user_input_as_number**2)
12 finally:
13    # This will be executed whether or not an exception is raised.
14    print('Thank you')
```

When using a `try` clause, you should have as little as possible in the `try` block. If too many things happen in that block, you risk handling an unexpected exception.

If the `try` block can fail if various way, you can handle different exceptions in the same `try` clause:

It is also possible not to specify a particular exception in the `except` statement. In this case, any exception will be handled. Such bare `except` statement should be avoided, though, as they can easily mask bugs.

## E.2 Raising our own exceptions

Can our program deliberately cause its own exceptions? If our program detects an error condition, we can **raise** an exception. Here is an example that gets input from the user and checks that the number is non-negative:

```
1 def get_age():
2     age = int(input("Please enter your age: "))
3     if age < 0:
4         # Create a new instance of an exception
5         my_error = ValueError("{0} is not a valid age".format(age))
6         raise my_error
7     return age
```

Line 5 creates an exception object, in this case, a `ValueError` object, which encapsulates specific information about the error. Assume that in this case function A called B which called C which called D which called `get_age`. The `raise` statement on line 6 carries this object out as a kind of “return value”, and immediately exits from `get_age()` to its caller D. Then D again exits to its caller C, and C exits to B and so on, each returning the exception object to their caller, until it encounters a `try ... except` that can handle the exception. We call this “unwinding the call stack”.

`ValueError` is one of the built-in exception types which most closely matches the kind of error we want to raise. The complete listing of built-in exceptions can be found at the [Built-in Exceptions](#) section of the [Python Library Reference](#), again by Python’s creator, Guido van Rossum.

If the function that called `get_age` (or its caller, or their caller, ...) handles the error, then the program can carry on running; otherwise, Python prints the traceback and exits:

```
>>> get_age()
Please enter your age: 42
42
>>> get_age()
Please enter your age: -2
Traceback (most recent call last):
  File "<interactive input>", line 1, in <module>
  File "learn_exceptions.py", line 4, in get_age
    raise ValueError("{0} is not a valid age".format(age))
ValueError: -2 is not a valid age
```

The error message includes the exception type and the additional information that was provided when the exception object was first created.

It is often the case that lines 5 and 6 (creating the exception object, then raising the exception) are combined into a single statement, but there are really two different and independent things happening, so perhaps it makes sense to keep the two steps separate when we first learn to work with exceptions. Here we show it all in a single statement:

```
1 raise ValueError("{0} is not a valid age".format(age))
```

## E.3 Revisiting an earlier example

Using exception handling, we can now modify our `recursion_depth` example from the previous chapter so that it stops when it reaches the maximum recursion depth allowed:

```
1 def recursion_depth(number):
2     print("Recursion depth number", number)
3     try:
4         recursion_depth(number + 1)
5     except:
6         print("I cannot go any deeper into this wormhole.")
7
8 recursion_depth(0)
```

Run this version and observe the results.

## E.4 The `finally` clause of the `try` statement

A common programming pattern is to grab a resource of some kind — e.g. we create a window for turtles to draw on, or we dial up a connection to our internet service provider, or we may open a file for writing. Then we perform some

computation which may raise an exception, or may work without any problems.

Whatever happens, we want to “clean up” the resources we grabbed — e.g. close the window, disconnect our dial-up connection, or close the file. The `finally` clause of the `try` statement is the way to do just this. Consider this (somewhat contrived) example:

```
1  import turtle
2  import time
3
4  def show_poly():
5      try:
6          win = turtle.Screen()    # Grab/create a resource, e.g. a window
7          tess = turtle.Turtle()
8
9          # This dialog could be cancelled,
10         # or the conversion to int might fail, or n might be zero.
11         n = int(input("How many sides do you want in your polygon?"))
12         angle = 360 / n
13         for i in range(n):        # Draw the polygon
14             tess.forward(10)
15             tess.left(angle)
16         time.sleep(3)              # Make program wait a few seconds
17     finally:
18         win.bye()                  # Close the turtle's window
19
20 show_poly()
21 show_poly()
22 show_poly()
```

In lines 20–22, `show_poly` is called three times. Each one creates a new window for its turtle, and draws a polygon with the number of sides input by the user. But what if the user enters a string that cannot be converted to an `int`? What if they close the dialog? We’ll get an exception, *but even though we’ve had an exception, we still want to close the turtle’s window*. Lines 17–18 does this for us. Whether we complete the statements in the `try` clause successfully or not, the `finally` block will always be executed.

Notice that the exception is still unhandled — only an `except` clause can handle an exception, so our program will still crash. But at least its turtle window will be closed before it crashes!

## E.5 Glossary

**exception** An error that occurs at runtime.

**handle an exception** To prevent an exception from causing our program to crash, by wrapping the block of code in a `try... except` construct.

**raise** To create a deliberate exception by using the `raise` statement.

## E.6 Exercises

1. Write a function named `readposint` that uses the `input` dialog to prompt the user for a positive integer and then checks the input to confirm that it meets the requirements. It should be able to handle inputs that cannot be converted to `int`, as well as negative `ints`, and edge cases (e.g. when the user closes the dialog, or does not enter anything at all.)

## APPENDIX F

---

### Fitting

---

Suppose we want to determine the gravitational acceleration. To this end, we could drop an object from the building and measure how long it takes for the object to reach the ground with a stopwatch. Newton's laws predict the following model:

$$h = \frac{1}{2}gt^2$$

where  $h$  is the height from which we dropped the object, and  $t$  the time it takes to hit the ground. So from our one measurement we could now calculate the gravitational acceleration  $g$ .

We can measure the height very accurately, but since we use a stopwatch to measure time, that value is a lot less reliable because you might have started and stopped your stopwatch at the wrong moments. Therefore, the result will not be very accurate. To make the value more accurate we should repeat the same measurement  $n$  times to obtain an average  $t$  and use that instead. We'll get back to this later.

For now we are more interested in the question: is this model correct? To test this question, we drop our object from different heights, doing multiple measurements for each height to get reliable values. The data, obtained by simulation for health and safety reasons, are given in the following table:

y	t	n
10	1.4	5
20	2.1	3
30	2.6	8
40	3.0	15
50	3.3	30

Since the model predicts a parabola, we want to fit the data to this model to see how good it works. It might be a bit confusing, but  $h$  is our x axis, and  $t$  is the y axis.

We use the `symfit` package to do our fitting. You can find the installation instructions [here](#).

To fit the data to the model we run the following code:



```
import numpy as np
from sympfit import Variable, Parameter, Fit, Model, sqrt

t_data = np.array([1.4, 2.1, 2.6, 3.0, 3.3])
h_data = np.array([10, 20, 30, 40, 50])

# We now define our model
h = Variable('h')
t = Variable('t')
g = Parameter('g')

t_model = Model({t: sqrt(2 * h / g)})

fit = Fit(t_model, h=h_data, t=t_data)
fit_result = fit.execute()
print(fit_result)
```

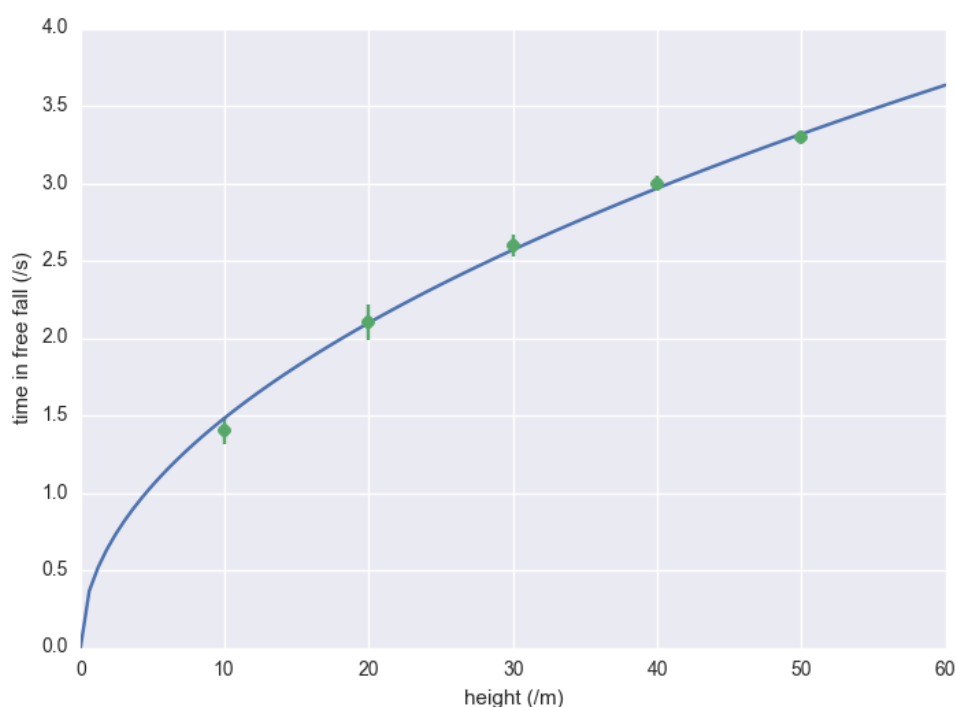
Looking at these results, we see that  $g = 9.09 \pm 0.15$  for this dataset. In order to plot this result alongside the data, we need to calculate values for the model. In the same script, we can do:

```
# Make an array from 0 to 50 in 1000 steps
h_range = np.linspace(0, 50, 1000)
fit_data = t_model(h=h_range, g=fit_result.value(g))
t_fit = fit_data.t
```

---

**Note:** When calling a `sympfit.Model`, a `namedtuple` is returned with all the components of the model evaluated. In this case there is only one component, and can be accessed by requesting `fit_data.t` or `fit_data[0]`.

---



This gives the model evaluated at all the points in `h_range`. Making the actual plot is left to you as an exercise. We see that we can request the fitted value of `g` by calling `fit_result.value(g)`, which returns 9.09.

Let's think for a second about the implications. The value of  $g$  is  $g = 9.81$  in the Netherlands. Based on our result, the textbooks should be rewritten because that value is extremely unlikely to be true given the small standard deviation in our data. It is at this point that we remember that our data point were not infinitely precise: we took many measurements and averaged them. This means there is an uncertainty in each of our data points. We will now account for this additional uncertainty and see what this does to our conclusion. To do this we first have to describe how the fitting actually works.

## F.1 How does it work?

In fitting we want to find values for the parameters such that the differences between the model and the data are as small as possible. The differences (residuals) are easy to calculate:

$$f(x_i, \vec{p}) - y_i$$

Here we have written the parameters as a vector  $\vec{p}$ , to indicate that we can have multiple parameters.  $x_i$  and  $y_i$  are the  $x$  and  $y$  coordinates of the  $i$ 'th datapoint. However, if we were to minimize the sum over all these differences we would have a problem, because these differences can be either positive or negative. This means there's many ways to add these values and get zero out of the sum. We therefore take the sum over the residuals squared:

$$Q^2 = \sum_{i=1}^n (f(x_i, \vec{p}) - y_i)^2$$

Now if we minimize  $Q^2$ , we get the best possible values for our parameters. The fitting algorithm actually just takes some values for the parameters, calculates  $Q^2$ , then changes the values slightly by adding or subtracting a number, and checks if this new value is smaller than the old one. If this is true it keeps going in the same direction until the value of  $Q^2$  starts to increase. That's when you know you've hit a minimum. Of course the trick is to do this smartly, and a lot of algorithms have been developed in order to do this.

## F.2 Propagating Uncertainties

In the example above we the fitting process assumed that every measurement was equally reliable. But this is not true. By repeating a measurement and averaging the result, we can improve the accuracy. So in our example, we dropped our object from every height a couple of times and took the average. Therefore, we want to assign a weight depending on how accurate the average value for that height is. Statistically the weight  $w_i$  to use is  $w_i = \frac{1}{\sigma_i^2}$ , where  $\sigma_i$  is the standard deviation for each point.

Our sum to minimize now changes to:

$$\chi^2 = \sum_{i=1}^n w_i (f(x_i, \vec{p}) - y_i)^2 = \sum_{i=1}^n \frac{(f(x_i, \vec{p}) - y_i)^2}{\sigma_i^2}$$

But how do we know the standard deviation in the mean value we calculate for every height? Suppose the standard deviation of our stopwatch is  $\sigma_{stopwatch} = 0.2$ . If we do  $n$  measurements from the same height, the average time is found by calculating

$$\bar{t} = \frac{1}{n} \sum_{i=1}^n t_i$$

It can be shown that the standard deviation of the mean is now:

$$\sigma_{\bar{t}} = \frac{\sigma_{stopwatch}}{\sqrt{n}}$$

So we see that by increasing the amount of measurements, we can decrease the uncertainty in  $\bar{t}$ . Our simulated data now changes to:

y	t	n	$\sigma_t$
10	1.4	5	0.089
20	2.1	3	0.115
30	2.6	8	0.071
40	3.0	15	0.052
50	3.3	30	0.037

The values of  $\sigma_t$  have been calculated by using the above formula. Let's fit to this new data set using *symfit*. Notice that there are some small differences to the code:

```
import numpy as np
from symfit import Variable, Parameter, Fit, Model, sqrt

t_data = np.array([1.4, 2.1, 2.6, 3.0, 3.3])
h_data = np.array([10, 20, 30, 40, 50])
n = np.array([5, 3, 8, 15, 30])
sigma = 0.2
sigma_t = sigma / np.sqrt(n)

# We now define our model
h = Variable('h')
t = Variable('t')
g = Parameter('g')
t_model = Model({t: sqrt(2 * h / g)})

fit = Fit(t_model, h=h_data, t=t_data, sigma_t=sigma_t)
fit_result = fit.execute()
print(fit_result)
```

---

**Note:** Looking at the initiation of `Fit`, we see that uncertainties can be provided to `Variable`'s by prepending their name with `sigma_`, in this case `sigma_t`.

---

Including these uncertainties in the fit yields  $g = 9.10 \pm 0.16$ . The accepted value of  $g = 9.81$  is well outside the uncertainty in this data. Therefore the textbooks must be rewritten!

This example shows the importance of propagating your errors consistently. (And of the importance of performing the actual measurement as the author of a chapter on error propagation so you don't end up claiming the textbooks have to be rewritten.)

## F.3 More on `symfit`

There are a lot more features in `symfit` to help you on your quest to fitting the universe. You can find the tutorial [here](#).

It is recommended you read this as well before starting to fit your own data.

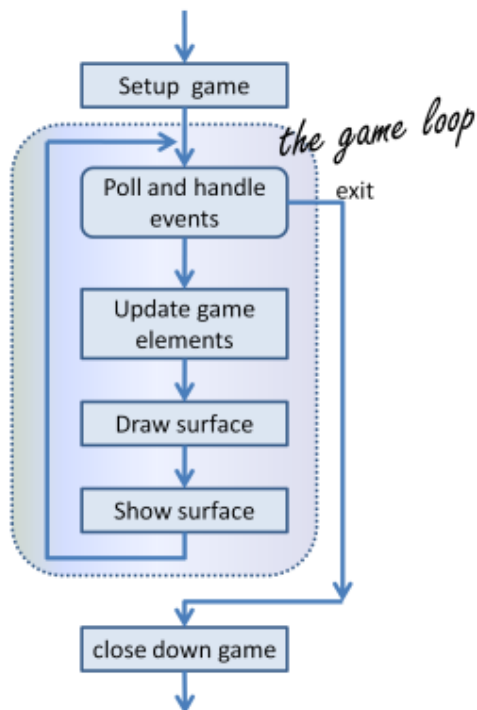


PyGame is a package that is not part of the standard Python distribution, so if you do not already have it installed (i.e. `import pygame` fails), download and install a suitable version from <http://pygame.org/download.shtml>. These notes are based on PyGame 1.9.1, the most recent version at the time of writing.

PyGame comes with a substantial set of tutorials, examples, and help, so there is ample opportunity to stretch yourself on the code. You may need to look around a bit to find these resources, though: if you've installed PyGame on a Windows machine, for example, they'll end up in a folder like `C:\Python31\Lib\site-packages\pygame\` where you will find directories for *docs* and *examples*.

### G.1 The game loop

The structure of the games we'll consider always follows this fixed pattern:



In every game, in the *setup* section we'll create a window, load and prepare some content, and then enter the **game loop**. The game loop continuously does four main things:

- it **polls** for events — i.e. asks the system whether events have occurred — and responds appropriately,
- it updates whatever internal data structures or objects need changing,
- it draws the current state of the game into a (non-visible) surface,
- it puts the just-drawn surface on display.

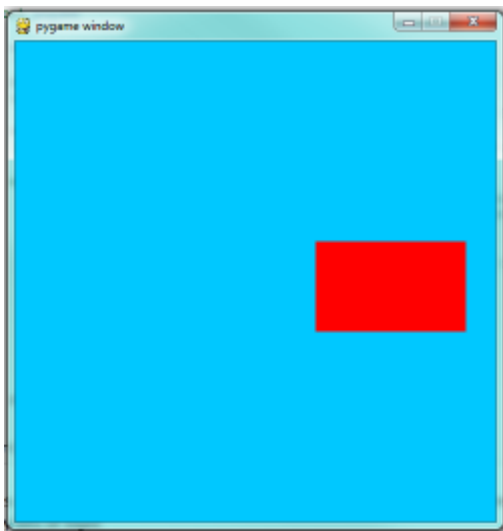
```
1  import pygame
2
3  def main():
4      """ Set up the game and run the main game loop """
5      pygame.init()          # Prepare the pygame module for use
6      surface_size = 480     # Desired physical surface size, in pixels.
7
8      # Create surface of (width, height), and its window.
9      main_surface = pygame.display.set_mode((surface_size, surface_size))
10
11     # Set up some data to describe a small rectangle and its color
12     small_rect = (300, 200, 150, 90)
13     some_color = (255, 0, 0)    # A color is a mix of (Red, Green, Blue)
14
15     while True:
16         event = pygame.event.poll()    # Look for any event
17         if event.type == pygame.QUIT:  # Window close button clicked?
18             break                      # ... leave game loop
19
20         # Update your game objects and data structures here...
21
22         # We draw everything from scratch on each frame.
23         # So first fill everything with the background color
```

(continues on next page)

(continued from previous page)

```
24     main_surface.fill((0, 200, 255))
25
26     # Overpaint a smaller rectangle on the main surface
27     main_surface.fill(some_color, small_rect)
28
29     # Now the surface is ready, tell pygame to display it!
30     pygame.display.flip()
31
32     pygame.quit()      # Once we leave the loop, close the window.
33
34 main()
```

This program pops up a window which stays there until we close it:



PyGame does all its drawing onto rectangular *surfaces*. After initializing PyGame at line 5, we create a window holding our main surface. The main loop of the game extends from line 15 to 30, with the following key bits of logic:

- First (line 16) we poll to fetch the next event that might be ready for us. This step will always be followed by some conditional statements that will determine whether any event that we’re interested in has happened. Polling for the event consumes it, as far as PyGame is concerned, so we only get one chance to fetch and use each event. On line 17 we test whether the type of the event is the predefined constant called `pygame.QUIT`. This is the event that we’ll see when the user clicks the close button on the PyGame window. In response to this event, we leave the loop.
- Once we’ve left the loop, the code at line 32 closes window, and we’ll return from function `main`. Your program could go on to do other things, or reinitialize pygame and create another window, but it will usually just end too.
- There are different kinds of events — key presses, mouse motion, mouse clicks, joystick movement, and so on. It is usual that we test and handle all these cases with new code squeezed in before line 19. The general idea is “handle events first, then worry about the other stuff”.
- At line 20 we’d update objects or data — for example, if we wanted to vary the color, position, or size of the rectangle we’re about to draw, we’d re-assign `some_color`, and `small_rect` here.
- A modern way to write games (now that we have fast computers and fast graphics cards) is to redraw everything from scratch on every iteration of the game loop. So the first thing we do at line 24 is fill the entire surface with a background color. The `fill` method of a surface takes two arguments — the color to use for filling, and the rectangle to be filled. But the second argument is optional, and if it is left out the entire surface is filled.



- In line 27 we fill a second rectangle, this time using `some_color`. The placement and size of the rectangle are given by the tuple `small_rect`, a 4-element tuple `(x, y, width, height)`.
- It is important to understand that the origin of the PyGame's surface is at the top left corner (unlike the turtle module that puts its origin in the middle of the screen). So, if you wanted the rectangle closer to the top of the window, you need to make its y coordinate smaller.
- If your graphics display hardware tries to read from memory at the same time as the program is writing to that memory, they will interfere with each other, causing video noise and flicker. To get around this, PyGame keeps two buffers in the main surface — the *back buffer* that the program draws to, while the *front buffer* is being shown to the user. Each time the program has fully prepared its back buffer, it flips the back/front role of the two buffers. So the drawing on lines 24 and 27 does not change what is seen on the screen until we `flip` the buffers, on line 30.

## G.2 Displaying images and text

To draw an image on the main surface, we load the image, say a beach ball, into its own new surface. The main surface has a `blit` method that copies pixels from the beach ball surface into its own surface. When we call `blit`, we can specify where the beach ball should be placed on the main surface. The term **blit** is widely used in computer graphics, and means *to make a fast copy of pixels from one area of memory to another*.

So in the setup section, before we enter the game loop, we'd load the image, like this:

```
1 ball = pygame.image.load("ball.png")
```

and after line 28 in the program above, we'd add this code to display our image at position (100,120):

```
1 main_surface.blit(ball, (100, 120))
```

To display text, we need do do three things. Before we enter the game loop, we instantiate a `font` object:

```
1 # Instantiate 16 point Courier font to draw text.
2 my_font = pygame.font.SysFont("Courier", 16)
```

and after line 28, again, we use the font's `render` method to create a new surface containing the pixels of the drawn text, and then, as in the case for images, we `blit` our new surface onto the main surface. Notice that `render` takes two extra parameters — the second tells it whether to carefully smooth edges of the text while drawing (this process is called *anti-aliasing*), and the second is the color that we want the text to be. Here we've used `(0, 0, 0)` which is black:

```
1 the_text = my_font.render("Hello, world!", True, (0,0,0))
2 main_surface.blit(the_text, (10, 10))
```

We'll demonstrate these two new features by counting the frames — the iterations of the game loop — and keeping some timing information. On each frame, we'll display the frame count, and the frame rate. We will only update the frame rate after every 500 frames, when we'll look at the timing interval and can do the calculations.

```
1 import pygame
2 import time
3
4 def main():
5
6     pygame.init()      # Prepare the PyGame module for use
7     main_surface = pygame.display.set_mode((480, 240))
8
```

(continues on next page)

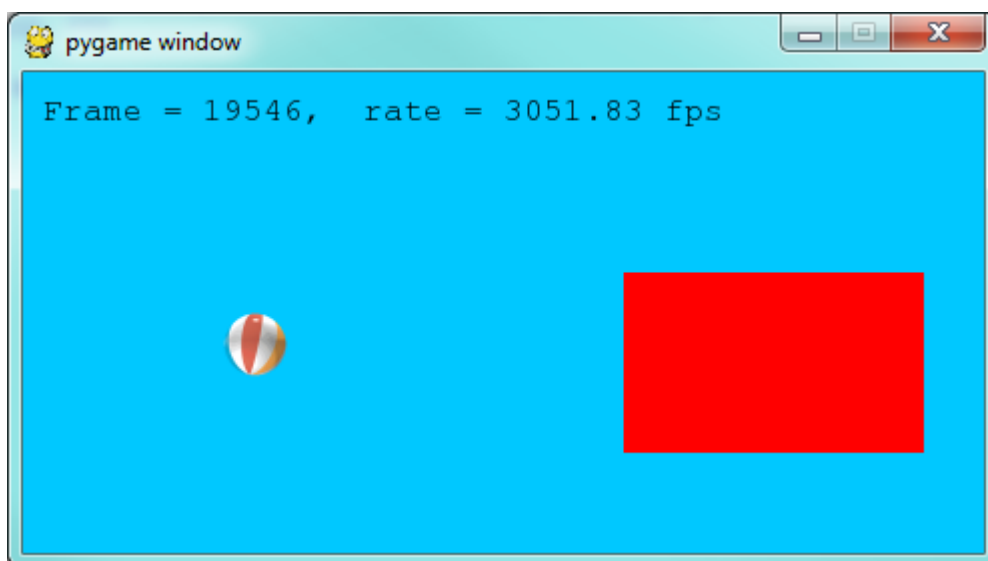
(continued from previous page)

```

9      # Load an image to draw. Substitute your own.
10     # PyGame handles gif, jpg, png, etc. image types.
11     ball = pygame.image.load("ball.png")
12
13     # Create a font for rendering text
14     my_font = pygame.font.SysFont("Courier", 16)
15
16     frame_count = 0
17     frame_rate = 0
18     t0 = time.clock()
19
20     while True:
21
22         # Look for an event from keyboard, mouse, joystick, etc.
23         ev = pygame.event.poll()
24         if ev.type == pygame.QUIT:    # Window close button clicked?
25             break                    # Leave game loop
26
27         # Do other bits of logic for the game here
28         frame_count += 1
29         if frame_count % 500 == 0:
30             t1 = time.clock()
31             frame_rate = 500 / (t1-t0)
32             t0 = t1
33
34         # Completely redraw the surface, starting with background
35         main_surface.fill((0, 200, 255))
36
37         # Put a red rectangle somewhere on the surface
38         main_surface.fill((255,0,0), (300, 100, 150, 90))
39
40         # Copy our image to the surface, at this (x,y) posn
41         main_surface.blit(ball, (100, 120))
42
43         # Make a new surface with an image of the text
44         the_text = my_font.render("Frame = {0}, rate = {1:.2f} fps"
45                                   .format(frame_count, frame_rate), True, (0,0,0))
46         # Copy the text surface to the main surface
47         main_surface.blit(the_text, (10, 10))
48
49         # Now that everything is drawn, put it on display!
50         pygame.display.flip()
51
52     pygame.quit()
53
54
55     main()

```

The frame rate is close to ridiculous — a lot faster than one’s eye can process frames. (Commercial video games usually plan their action for 60 frames per second (fps).) Of course, our rate will drop once we start doing something a little more strenuous inside our game loop.



## G.3 Drawing a board for the N queens puzzle

We previously solved our N queens puzzle. For the 8x8 board, one of the solutions was the list `[6, 4, 2, 0, 5, 7, 1, 3]`. Let's use that solution as testdata, and now use PyGame to draw that chessboard with its queens.

We'll create a new module for the drawing code, called `draw_queens.py`. When we have our test case(s) working, we can go back to our solver, import this new module, and add a call to our new function to draw a board each time a solution is discovered.

We begin with a background of black and red squares for the board. Perhaps we could create an image that we could load and draw, but that approach would need different background images for different size boards. Just drawing our own red and black rectangles of the appropriate size sounds like much more fun!

```
1 def draw_board(the_board):
2     """ Draw a chess board with queens, from the_board. """
3
4     pygame.init()
5     colors = [(255,0,0), (0,0,0)]    # Set up colors [red, black]
6
7     n = len(the_board)               # This is an NxN chess board.
8     surface_size = 480               # Proposed physical surface size.
9     square_size = surface_size // n  # sq_sz is length of a square.
10    surface_size = n * square_size    # Adjust to exactly fit n squares.
11
12    # Create the surface of (width, height), and its window.
13    surface = pygame.display.set_mode((surface_size, surface_size))
```

Here we precompute `square_size`, the integer size that each square will be, so that we can fit the squares nicely into the available window. So if we'd like the board to be 480x480, and we're drawing an 8x8 chessboard, then each square will need to have a size of 60 units. But we notice that a 7x7 board cannot fit nicely into 480 — we're going to get some ugly border that our squares don't fill exactly. So we recompute the surface size to exactly fit our squares before we create the window.

Now let's draw the squares, in the game loop. We'll need a nested loop: the outer loop will run over the rows of the chessboard, the inner loop over the columns:

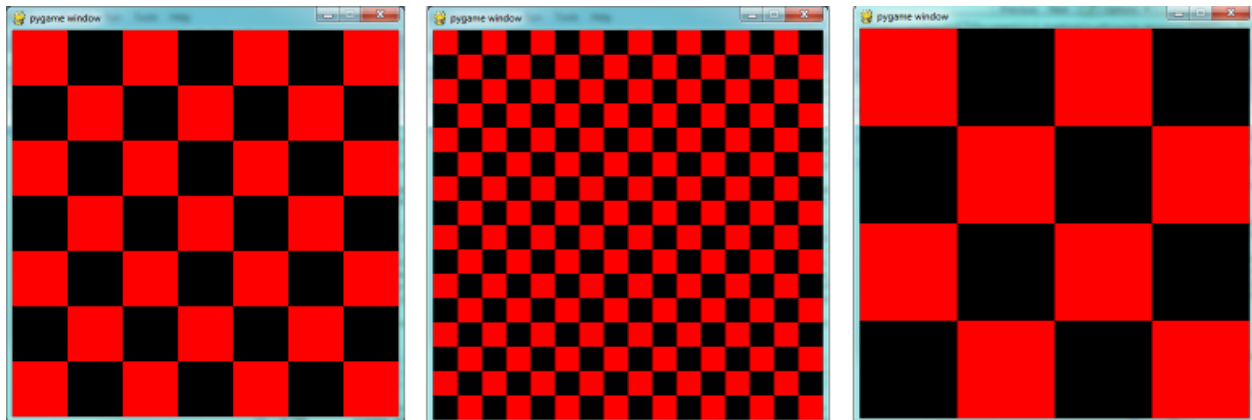
```

1  # Draw a fresh background (a blank chess board)
2  for row in range(n):          # Draw each row of the board.
3      color_index = row % 2      # Change starting color on each row
4      for col in range(n):      # Run through cols drawing squares
5          the_square = (col*square_size, row*square_size, square_
6      ↪size)
7          surface.fill(colors[color_index], the_square)
8          # now flip the color index for the next square
          c_index = (c_index + 1) % 2

```

There are two important ideas in this code: firstly, we compute the rectangle to be filled from the `row` and `col` loop variables, multiplying them by the size of the square to get their position. And, of course, each square is a fixed width and height. So `the_square` represents the rectangle to be filled on the current iteration of the loop. The second idea is that we have to alternate colors on every square. In the earlier setup code we created a list containing two colors, here we manipulate `color_index` (which will always either have the value 0 or 1) to start each row on a color that is different from the previous row's starting color, and to switch colors each time a square is filled.

This (together with the other fragments not shown to flip the surface onto the display) leads to the pleasing backgrounds like this, for different size boards:



Now, on to drawing the queens! Recall that our solution `[6, 4, 2, 0, 5, 7, 1, 3]` means that in column 0 of the board we want a queen at row 6, at column 1 we want a queen at row 4, and so on. So we need a loop running over each queen:

```

1  for (col, row) in enumerate(the_board):
2      # draw a queen at col, row...

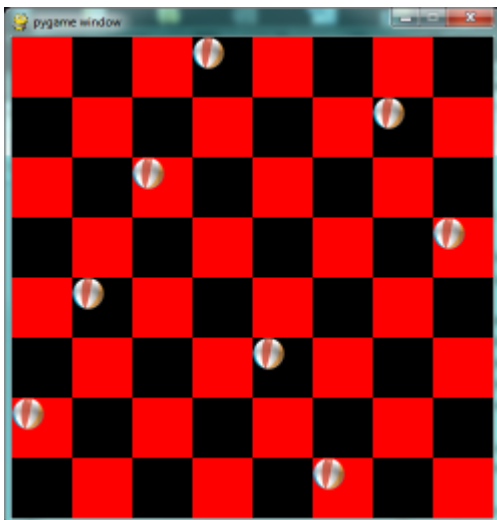
```

In this chapter we already have a beach ball image, so we'll use that for our queens. In the setup code before our game loop, we load the ball image (as we did before), and in the body of the loop, we add the line:

```

1  surface.blit(ball, (col * square_size, row * square_size))

```



We're getting there, but those queens need to be centred in their squares! Our problem arises from the fact that both the ball and the rectangle have their upper left corner as their reference points. If we're going to centre this ball in the square, we need to give it an extra offset in both the x and y direction. (Since the ball is round and the square is square, the offset in the two directions will be the same, so we'll just compute a single offset value, and use it in both directions.)

The offset we need is half the (size of the square less the size of the ball). So we'll precompute this in the game's setup section, after we've loaded the ball and determined the square size:

```
1 ball_offset = (square_size - ball.get_width()) // 2
```

Now we touch up the drawing code for the ball and we're done:

```
1 surface.blit(ball, (col * square_size + ball_offset, row * square_size +  
↪ball_offset))
```

We might just want to think about what would happen if the ball was bigger than the square. In that case, `ball_offset` would become negative. So it would still be centered in the square - it would just spill over the boundaries, or perhaps obscure the square entirely!

Here is the complete program:

```
1 import pygame
2
3 def draw_board(the_board):
4     """ Draw a chess board with queens, as determined by the the_board. """
5
6     pygame.init()
7     colors = [(255,0,0), (0,0,0)]    # Set up colors [red, black]
8
9     n = len(the_board)               # This is an NxN chess board.
10    surface_size = 480                # Proposed physical surface size.
11    square_size = surface_size // n   # sq_sz is length of a square.
12    surface_size = n * square_size    # Adjust to exactly fit n squares.
13
14    # Create the surface of (width, height), and its window.
15    surface = pygame.display.set_mode((surface_size, surface_size))
16
17    ball = pygame.image.load("ball.png")
18
```

(continues on next page)

(continued from previous page)

```

19     # Use an extra offset to centre the ball in its square.
20     # If the square is too small, offset becomes negative,
21     # but it will still be centered :-)
22     ball_offset = (square_size-ball.get_width()) // 2
23
24     while True:
25
26         # Look for an event from keyboard, mouse, etc.
27         event = pygame.event.poll()
28         if event.type == pygame.QUIT:
29             break;
30
31         # Draw a fresh background (a blank chess board)
32         for row in range(n):          # Draw each row of the board.
33             color_index = row % 2     # Alternate starting color
34             for col in range(n):      # Run through cols drawing squares
35                 the_square = (col*square_size, row*square_size, square_size,
↪square_size)
36                 surface.fill(colors[color_index], the_square)
37                 # Now flip the color index for the next square
38                 color_index = (color_index + 1) % 2
39
40         # Now that squares are drawn, draw the queens.
41         for (col, row) in enumerate(the_board):
42             surface.blit(ball,
43                           (col*square_size+ball_offset, row*square_size+ball_offset))
44
45         pygame.display.flip()
46
47
48     pygame.quit()
49
50 if __name__ == "__main__":
51     draw_board([0, 5, 3, 1, 6, 4, 2])    # 7 x 7 to test window size
52     draw_board([6, 4, 2, 0, 5, 7, 1, 3])
53     draw_board([9, 6, 0, 3, 10, 7, 2, 4, 12, 8, 11, 5, 1]) # 13 x 13
54     draw_board([11, 4, 8, 12, 2, 7, 3, 15, 0, 14, 10, 6, 13, 1, 5, 9])

```

There is one more thing worth reviewing here. The conditional statement on line 50 tests whether the name of the currently executing program is `__main__`. This allows us to distinguish whether this module is being run as a main program, or whether it has been imported elsewhere, and used as a module. If we run this module in Python, the test cases in lines 51-54 will be executed. However, if we import this module into another program (i.e. our N queens solver from earlier) the condition at line 50 will be false, and the statements on lines 51-54 won't run.

Previously, our main program looked like this:

```

1 def main():
2
3     board = list(range(8))    # Generate the initial permutation
4     num_found = 0
5     tries = 0
6     while num_found < 10:
7         random.shuffle(board)
8         tries += 1
9         if not has_clashes(board):
10             print("Found solution {0} in {1} tries.".format(board, tries))
11             tries = 0

```

(continues on next page)

(continued from previous page)

```
12         num_found += 1
13
14     main()
```

Now we just need two changes. At the top of that program, we import the module that we’ve been working on here (assume we called it `draw_queens`). (You’ll have to ensure that the two modules are saved in the same folder.) Then after line 10 here we add a call to draw the solution that we’ve just discovered:

```
draw_queens.draw_board(board)
```

And that gives a very satisfying combination of program that can search for solutions to the N queens problem, and when it finds each, it pops up the board showing the solution.

## G.4 Sprites

A sprite is an object that can move about in a game, and has internal behaviour and state of its own. For example, a spaceship would be a sprite, the player would be a sprite, and bullets and bombs would all be sprites.

Object oriented programming (OOP) is ideally suited to a situation like this: each object can have its own attributes and internal state, and a couple of methods. Let’s have some fun with our N queens board. Instead of placing the queen in her final position, we’d like to drop her in from the top of the board, and let her fall into position, perhaps bouncing along the way.

The first encapsulation we need is to turn each of our queens into an object. We’ll keep a list of all the active sprites (i.e. a list of queen objects), and arrange two new things in our game loop:

- After handling events, but before drawing, call an `update` method on every sprite. This will give each sprite a chance to modify its internal state in some way — perhaps change its image, or change its position, or rotate itself, or make itself grow a bit bigger or a bit smaller.
- Once all the sprites have updated themselves, the game loop can begin drawing - first the background, and then call a `draw` method on each sprite in turn, and delegate (hand off) the task of drawing to the object itself. This is in line with the OOP idea that we don’t say “Hey, *draw*, show this queen!”, but we prefer to say “Hey, *queen*, draw yourself!”.

We start with a simple object, no movement or animation yet, just scaffolding, to see how to fit all the pieces together:

```
1 class QueenSprite:
2
3     def __init__(self, img, target_posn):
4         """ Create and initialize a queen for this
5             target position on the board
6         """
7         self.image = img
8         self.target_posn = target_posn
9         self.position = target_posn
10
11     def update(self):
12         return # Do nothing for the moment.
13
14     def draw(self, target_surface):
15         target_surface.blit(self.image, self.position)
```

We’ve given the sprite three attributes: an image to be drawn, a target position, and a current position. If we’re going to move the spite about, the current position may need to be different from the target, which is where we want the queen

finally to end up. In this code at this time we've done nothing in the `update` method, and our `draw` method (which can probably remain this simple in future) simply draws itself at its current position on the surface that is provided by the caller.

With its class definition in place, we now instantiate our N queens, put them into a list of sprites, and arrange for the game loop to call the `update` and `draw` methods on each frame. The new bits of code, and the revised game loop look like this:

```

1  all_sprites = []           # Keep a list of all sprites in the game
2
3  # Create a sprite object for each queen, and populate our list.
4  for (col, row) in enumerate(the_board):
5      a_queen = QueenSprite(ball,
6                             (col*square_size+ball_offset, row*square_size+ball_
7                             ↳offset))
8      all_sprites.append(a_queen)
9
10 while True:
11     # Look for an event from keyboard, mouse, etc.
12     event = pygame.event.poll()
13     if event.type == pygame.QUIT:
14         break;
15
16     # Ask every sprite to update itself.
17     for sprite in all_sprites:
18         sprite.update()
19
20     # Draw a fresh background (a blank chess board)
21     # ... same as before ...
22
23     # Ask every sprite to draw itself.
24     for sprite in all_sprites:
25         sprite.draw(surface)
26
27     pygame.display.flip()
```

This works just like it did before, but our extra work in making objects for the queens has prepared the way for some more ambitious extensions.

Let us begin with a falling queen object. At any instant, it will have a velocity i.e. a speed, in a certain direction. (We are only working with movement in the y direction, but use your imagination!) So in the object's `update` method, we want to change its current position by its velocity. If our N queens board is floating in space, velocity would stay constant, but hey, here on Earth we have gravity! Gravity changes the velocity on each time interval, so we'll want a ball that speeds up as it falls further. Gravity will be constant for all queens, so we won't keep it in the instances — we'll just make it a variable in our module. We'll make one other change too: we will start every queen at the top of the board, so that it can fall towards its target position. With these changes, we now get the following:

```

1  gravity = 0.0001
2
3  class QueenSprite:
4
5      def __init__(self, img, target_posn):
6          self.image = img
7          self.target_position = target_position
8          (x, y) = target_position
9          self.position = (x, 0)      # Start ball at top of its column
10         self.y_velocity = 0         # with zero initial velocity
11
```

(continues on next page)



(continued from previous page)

```
12 def update(self):
13     self.y_velocity += gravity          # Gravity changes velocity
14     (x, y) = self.position
15     new_y_pos = y + self.y_velocity    # Velocity moves the ball
16     self.position = (x, new_y_pos)     # to this new position.
17
18 def draw(self, target_surface):        # Same as before.
19     target_surface.blit(self.image, self.position)
```

Making these changes gives us a new chessboard in which each queen starts at the top of its column, and speeds up, until it drops off the bottom of the board and disappears forever. A good start — we have movement!

The next step is to get the ball to bounce when it reaches its own target position. It is pretty easy to bounce something — you just change the sign of its velocity, and it will move at the same speed in the opposite direction. Of course, if it is travelling up towards the top of the board it will be slowed down by gravity. (Gravity always sucks down!) And you'll find it bounces all the way up to where it began from, reaches zero velocity, and starts falling all over again. So we'll have bouncing balls that never settle.

A realistic way to settle the object is to lose some energy (probably to friction) each time it bounces — so instead of simply reversing the sign of the velocity, we multiply it by some fractional factor — say -0.65. This means the ball only retains 65% of its energy on each bounce, so it will, as in real life, stop bouncing after a short while, and settle on its “ground”.

The only changes are in the update method, which now looks like this:

```
1 def update(self):
2     self.y_velocity += gravity
3     (x, y) = self.position
4     new_y_pos = y + self.y_velocity
5     (target_x, target_y) = self.target_posn    # Unpack the position
6     dist_to_go = target_y - new_y_pos          # How far to our floor?
7
8     if dist_to_go < 0:                          # Are we under floor?
9         self.y_velocity = -0.65 * self.y_velocity    # Bounce
10        new_y_pos = target_y + dist_to_go          # Move back above floor
11
12    self.position = (x, new_y_pos)              # Set our new position.
```

Heh, heh, heh! We're not going to show animated screenshots, so copy the code into your Python environment and see for yourself.

## G.5 Events

The only kind of event we're handled so far has been the QUIT event. But we can also detect keydown and keyup events, mouse motion, and mousebutton down or up events. Consult the PyGame documentation and follow the link to Event.

When your program polls for and receives an event object from PyGame, its event type will determine what secondary information is available. Each event object carries a *dictionary* (which you may only cover in due course in these notes). The dictionary holds certain *keys* and *values* that make sense for the type of event.

For example, if the type of event is MOUSEMOTION, we'll be able to find the mouse position and information about the state of the mouse buttons in the dictionary attached to the event. Similarly, if the event is KEYDOWN, we can learn from the dictionary which key went down, and whether any modifier keys (shift, control, alt, etc.) are also down. You also get events when the game window becomes active (i.e. gets focus) or loses focus.

The event object with type `NOEVENT` is returned if there are no events waiting. Events can be printed, allowing you to experiment and play around. So dropping these lines of code into the game loop directly after polling for any event is quite informative:

```
1  if event.type != pygame.NOEVENT:    # Only print if it is interesting!
2      print(event)
```

With this in place, hit the space bar and the escape key, and watch the events you get. Click your three mouse buttons. Move your mouse over the window. (This causes a vast cascade of events, so you may also need to filter those out of the printing.) You'll get output that looks something like this:

```
<Event(17-VideoExpose {})>
<Event(1-ActiveEvent {'state': 1, 'gain': 0})>
<Event(2-KeyDown {'scancode': 57, 'key': 32, 'unicode': ' ', 'mod': 0})>
<Event(3-KeyUp {'scancode': 57, 'key': 32, 'mod': 0})>
<Event(2-KeyDown {'scancode': 1, 'key': 27, 'unicode': '\x1b', 'mod': 0})>
<Event(3-KeyUp {'scancode': 1, 'key': 27, 'mod': 0})>
...
<Event(4-MouseMotion {'buttons': (0, 0, 0), 'pos': (323, 194), 'rel': (-3, -1)})>
<Event(4-MouseMotion {'buttons': (0, 0, 0), 'pos': (322, 193), 'rel': (-1, -1)})>
<Event(4-MouseMotion {'buttons': (0, 0, 0), 'pos': (321, 192), 'rel': (-1, -1)})>
<Event(4-MouseMotion {'buttons': (0, 0, 0), 'pos': (319, 192), 'rel': (-2, 0)})>
<Event(5-MouseButtonDown {'button': 1, 'pos': (319, 192)})>
<Event(6-MouseButtonUp {'button': 1, 'pos': (319, 192)})>
<Event(4-MouseMotion {'buttons': (0, 0, 0), 'pos': (319, 191), 'rel': (0, -1)})>
<Event(5-MouseButtonDown {'button': 2, 'pos': (319, 191)})>
<Event(5-MouseButtonDown {'button': 5, 'pos': (319, 191)})>
<Event(6-MouseButtonUp {'button': 5, 'pos': (319, 191)})>
<Event(6-MouseButtonUp {'button': 2, 'pos': (319, 191)})>
<Event(5-MouseButtonDown {'button': 3, 'pos': (319, 191)})>
<Event(6-MouseButtonUp {'button': 3, 'pos': (319, 191)})>
...
<Event(1-ActiveEvent {'state': 1, 'gain': 0})>
<Event(12-Quit {})>
```

So let us now make these changes to the code near the top of our game loop:

```
1  while True:
2
3      # Look for an event from keyboard, mouse, etc.
4      ev = pygame.event.poll()
5      if event.type == pygame.QUIT:
6          break;
7      if event.type == pygame.KEYDOWN:
8          key = ev.dict["key"]
9          if key == 27:                                # On Escape key ...
10             break                                     # leave the game loop.
11         if key == ord("r"):
12             colors[0] = (255, 0, 0)    # Change to red + black.
13         elif key == ord("g"):
14             colors[0] = (0, 255, 0)    # Change to green + black.
15         elif key == ord("b"):
16             colors[0] = (0, 0, 255)    # Change to blue + black.
17
18     if event.type == pygame.MOUSEBUTTONDOWN: # Mouse gone down?
19         posn_of_click = event.dict["pos"]    # Get the coordinates.
20         print(posn_of_click)                 # Just print them.
```

Lines 7-16 show typical processing for a KEYDOWN event — if a key has gone down, we test which key it is, and take some action. With this in place, we have another way to quit our queens program — by hitting the escape key. Also, we can use keys to change the color of the board that is drawn.

Finally, at line 20, we respond (pretty lamely) to the mouse button going down.

As a final exercise in this section, we'll write a better response handler to mouse clicks. What we will do is figure out if the user has clicked the mouse on one of our sprites. If there is a sprite under the mouse when the click occurs, we'll send the click to the sprite and let it respond in some sensible way.

We'll begin with some code that finds out which sprite is under the clicked position, perhaps none! We add a method to the class, `contains_point`, which returns True if the point is within the rectangle of the sprite:

```
1  def contains_point(self, point):
2      """ Return True if my sprite rectangle contains point pt """
3      (my_x, my_y) = self.position
4      my_width = self.image.get_width()
5      my_height = self.image.get_height()
6      (x, y) = point
7      return ( x >= my_x and x < my_x + my_width and
8              y >= my_y and y < my_y + my_height)
```

Now in the game loop, once we've seen the mouse event, we determine which queen, if any, should be told to respond to the event:

```
1  if ev.type == pygame.MOUSEBUTTONDOWN:
2      posn_of_click = event.dict["pos"]
3      for sprite in all_sprites:
4          if sprite.contains_point(posn_of_click):
5              sprite.handle_click()
6              break
```

And the final thing is to write a new method called `handle_click` in the `QueenSprite` class. When a sprite is clicked, we'll just add some velocity in the up direction, i.e. kick it back into the air.

```
1  def handle_click(self):
2      self.y_velocity += -0.3  # Kick it up
```

With these changes we have a playable game! See if you can keep all the balls on the move, not allowing any one to settle!

## G.6 A wave of animation

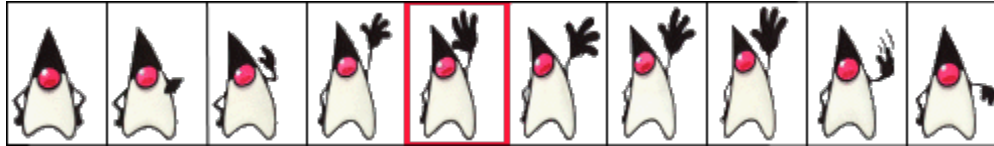
Many games have sprites that are animated: they crouch, jump and shoot. How do they do that?

Consider this sequence of 10 images: if we display them in quick succession, Duke will wave at us. (Duke is a friendly visitor from the kingdom of Javaland.)



A compound image containing smaller *patches* which are intended for animation is called a **sprite sheet**. Download this sprite sheet by right-clicking in your browser and saving it in your working directory with the name `duke_spritesheet.png`.

The sprite sheet has been quite carefully prepared: each of the 10 patches are spaced exactly 50 pixels apart. So, assuming we want to draw patch number 4 (numbering from 0), we want to draw only the rectangle that starts at x position 200, and is 50 pixels wide, within the sprite sheet. Here we’ve shown the patches and highlighted the patch we want to draw.



The `blit` method we’ve been using — for copying pixels from one surface to another — can copy a sub-rectangle of the source surface. So the grand idea here is that each time we draw Duke, we won’t blit the whole sprite sheet. Instead we’ll provide an extra rectangle argument that determines which portion of the sprite sheet will be blitted.

We’re going to add new code in this section to our existing N queens drawing game. What we want is to put some instances of Duke on the chessboard somewhere. If the user clicks on one of them, we’ll get him to respond by waving back, for one cycle of his animation.

But before we do that, we need another change. Up until now, our game loop has been running at really fast frame rates that are unpredictable. So we’ve chosen some *magic numbers* for gravity and for bouncing and kicking the ball on the basis of trial-and-error. If we’re going to start animating more sprites, we need to tame our game loop to operate at a fixed, known frame rate. This will allow us to plan our animation better.

PyGame gives us the tools to do this in just two lines of code. In the setup section of the game, we instantiate a new `Clock` object:

```
1 my_clock = pygame.time.Clock()
```

and right at the bottom of the game loop, we call a method on this object that limits the frame rate to whatever we specify. So let’s plan our game and animation for 60 frames per second, by adding this line at the bottom of our game loop:

```
1 my_clock.tick(60) # Waste time so that frame rate becomes 60 fps
```

You’ll find that you have to go back and adjust the numbers for gravity and kicking the ball now, to match this much slower frame rate. When we plan an animation so that it only works sensibly at a fixed frame rate, we say that we’ve *baked* the animation. In this case we’re baking our animations for 60 frames per second.

To fit into the existing framework that we already have for our queens board, we want to create a `DukeSprite` class that has all the same methods as the `QueenSprite` class. Then we can add one or more Duke instances onto our list of `all_sprites`, and our existing game loop will then call methods of the Duke instance. Let us start with skeleton scaffolding for the new class:

```
1 class DukeSprite:
2
3     def __init__(self, img, target_position):
4         self.image = img
5         self.position = target_position
6
7     def update(self):
8         return
9
10    def draw(self, target_surface):
11        return
12
13    def handle_click(self):
14        return
15
```

(continues on next page)

(continued from previous page)

```
16     def contains_point(self, pt):
17         # Use code from QueenSprite here
18         return
```

The only changes we'll need to the existing game are all in the setup section. We load up the new sprite sheet and instantiate a couple of instances of Duke, at the positions we want on the chessboard. So before entering the game loop, we add this code:

```
1  # Load the sprite sheet
2  duke_sprite_sheet = pygame.image.load("duke_spritesheet.png")
3
4  # Instantiate two duke instances, put them on the chessboard
5  duke1 = DukeSprite(duke_sprite_sheet, (square_size*2, 0))
6  duke2 = DukeSprite(duke_sprite_sheet, (square_size*5, sq_sz))
7
8  # Add them to the list of sprites which our game loop manages
9  all_sprites.append(duke1)
10 all_sprites.append(duke2)
```

Now the game loop will test if each instance has been clicked, will call the click handler for that instance. It will also call update and draw for all sprites. All the remaining changes we need to make will be made in the methods of the DukeSprite class.

Let's begin with drawing one of the patches. We'll introduce a new attribute `curr_patch_num` into the class. It holds a value between 0 and 9, and determines which patch to draw. So the job of the draw method is to compute the sub-rectangle of the patch to be drawn, and to blit only that portion of the spritesheet:

```
1  def draw(self, target_surface):
2      patch_rect = (self.curr_patch_num * 50, 0,
3                  50, self.image.get_height())
4      target_surface.blit(self.image, self.posn, patch_rect)
```

Now on to getting the animation to work. We need to arrange logic in update so that if we're busy animating, we change the `curr_patch_num` every so often, and we also decide when to bring Duke back to his rest position, and stop the animation. An important issue is that the game loop frame rate — in our case 60 fps — is not the same as the *animation rate* — the rate at which we want to change Duke's animation patches. So we'll plan Duke wave's animation cycle for a duration of 1 second. In other words, we want to play out Duke's 10 animation patches over 60 calls to update. (This is how the baking of the animation takes place!) So we'll keep another animation frame counter in the class, which will be zero when we're not animating, and each call to update will increment the counter up to 59, and then back to 0. We can then divide that animation counter by 6, to set the `curr_patch_num` variable to select the patch we want to show.

```
1  def update(self):
2      if self.anim_frame_count > 0:
3          self.anim_frame_count = (self.anim_frame_count + 1) % 60
4          self.curr_patch_num = self.anim_frame_count // 6
```

Notice that if `anim_frame_count` is zero, i.e. Duke is at rest, nothing happens here. But if we start the counter running, it will count up to 59 before settling back to zero. Notice also, that because `anim_frame_count` can only be a value between 0 and 59, the `curr_patch_num` will always stay between 0 and 9. Just what we require!

Now how do we trigger the animation, and start it running? On the mouse click.

```
1  def handle_click(self):
2      if self.anim_frame_count == 0:
3          self.anim_frame_count = 5
```

Two things of interest here. We only start the animation if Duke is at rest. Clicks on Duke while he is already waving get ignored. And when we do start the animation, we set the counter to 5 — this means that on the very next call to `update` the counter becomes 6, and the image changes. If we had set the counter to 1, we would have needed to wait for 5 more calls to `update` before anything happened — a slight lag, but enough to make things feel sluggish.

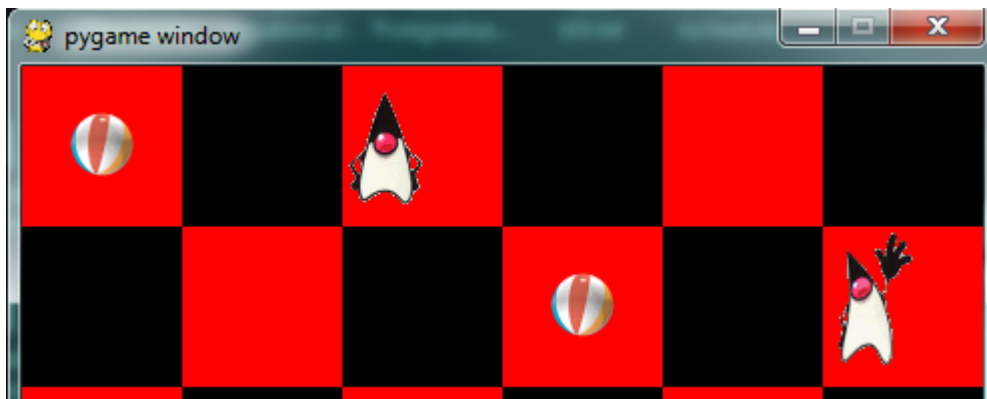
The final touch-up is to initialize our two new attributes when we instantiate the class. Here is the code for the whole class now:

```

1  class DukeSprite:
2
3      def __init__(self, img, target_posn):
4          self.image = img
5          self.position = target_posn
6          self.anim_frame_count = 0
7          self.curr_patch_num = 0
8
9      def update(self):
10         if self.anim_frame_count > 0:
11             self.anim_frame_count = (self.anim_frame_count + 1) % 60
12             self.curr_patch_num = self.anim_frame_count // 6
13
14         def draw(self, target_surface):
15             patch_rect = (self.curr_patch_num * 50, 0,
16                          50, self.image.get_height())
17             target_surface.blit(self.image, self.posn, patch_rect)
18
19         def contains_point(self, pt):
20             """ Return True if my sprite rectangle contains pt """
21             (my_x, my_y) = self.posn
22             my_width = self.image.get_width()
23             my_height = self.image.get_height()
24             (x, y) = pt
25             return ( x >= my_x and x < my_x + my_width and
26                     y >= my_y and y < my_y + my_height)
27
28         def handle_click(self):
29             if self.anim_frame_count == 0:
30                 self.anim_frame_count = 5

```

Now we have two extra Duke instances on our chessboard, and clicking on either causes that instance to wave.



## G.7 Aliens - a case study

Find the example games with the PyGame package, (On a windows system, something like C:\Python3\Lib\site-packages\pygame\examples) and play the Aliens game. Then read the code, in an editor or Python environment that shows line numbers.

It does a number of much more advanced things that we do, and relies on the PyGame framework for more of its logic. Here are some of the points to notice:

- The frame rate is deliberately constrained near the bottom of the game loop at line 311. If we change that number we can make the game very slow or unplayably fast!
- There are different kinds of sprites: Explosions, Shots, Bombs, Aliens and a Player. Some of these have more than one image — by swapping the images, we get animation of the sprites, i.e. the Alien spacecraft lights change, and this is done at line 112.
- Different kinds of objects are referenced in different groups of sprites, and PyGame helps maintain these. This lets the program check for collisions between, say, the list of shots fired by the player, and the list of spaceships that are attacking. PyGame does a lot of the hard work for us.
- Unlike our game, objects in the Aliens game have a limited lifetime, and have to get killed. For example, if we shoot, a Shot object is created — if it reaches the top of the screen without exploding against anything, it has to be removed from the game. Lines 141-142 do this. Similarly, when a falling bomb gets close to the ground (line 156), it instantiates a new Explosion sprite, and the bomb kills itself.
- There are random timings that add to the fun — when to spawn the next Alien, when an Alien drops the next bomb, etc.
- The game plays sounds too: a less-than-relaxing loop sound, plus sounds for the shots and explosions.

## G.8 Reflections

Object oriented programming is a good organizational tool for software. In the examples in this chapter, we've started to use (and hopefully appreciate) these benefits. Here we had N queens each with its own state, falling to its own floor level, bouncing, getting kicked, etc. We might have managed without the organizational power of objects — perhaps we could have kept lists of velocities for each queen, and lists of target positions, and so on — our code would likely have been much more complicated, ugly, and a lot poorer!

## G.9 Glossary

**animation rate** The rate at which we play back successive patches to create the illusion of movement. In the sample we considered in this chapter, we played Duke's 10 patches over the duration of one second. Not the same as the frame rate.

**baked animation** An animation that is designed to look good at a predetermined fixed frame rate. This reduces the amount of computation that needs to be done when the game is running. High-end commercial games usually bake their animations.

**blit** A verb used in computer graphics, meaning to make a fast copy of an image or pixels from a sub-rectangle of one image or surface to another surface or image.

**frame rate** The rate at which the game loop executes and updates the display.

**game loop** A loop that drives the logic of a game. It will usually poll for events, then update each of the objects in the game, then get everything drawn, and then put the newly drawn frame on display.

**pixel** A single picture element, or dot, from which images are made.

**poll** To ask whether something like a keypress or mouse movement has happened. Game loops usually poll to discover what events have occurred. This is different from event-driven programs like the ones seen in the chapter titled “Events”. In those cases, the button click or keypress event triggers the call of a handler function in your program, but this happens behind your back.

**sprite** An active agent or element in a game, with its own state, position and behaviour.

**surface** This is PyGame’s term for what the Turtle module calls a *canvas*. A surface is a rectangle of pixels used for displaying shapes and images.

## G.10 Exercises

1. Have fun with Python, and with PyGame.
2. We deliberately left a bug in the code for animating Duke. If you click on one of the chessboard squares to the right of Duke, he waves anyway. Why? Find a one-line fix for the bug.
3. Use your preferred search engine to search their image library for “sprite sheet playing cards”. Create a list [0..51] to represent an encoding of the 52 cards in a deck. Shuffle the cards, slice off the top five as your hand in a poker deal. Display the hand you have been dealt.
4. So the Aliens game is in outer space, without gravity. Shots fly away forever, and bombs don’t speed up when they fall. Add some gravity to the game. Decide if you’re going to allow your own shots to fall back on your head and kill you.
5. Those pesky Aliens seem to pass right through each other! Change the game so that they collide, and destroy each other in a mighty explosion.





---

## Plotting data with matplotlib

---

### H.1 Introduction

There are many scientific plotting packages. In this chapter we focus on **matplotlib**, chosen because it is the *de facto* plotting library and integrates very well with Python.

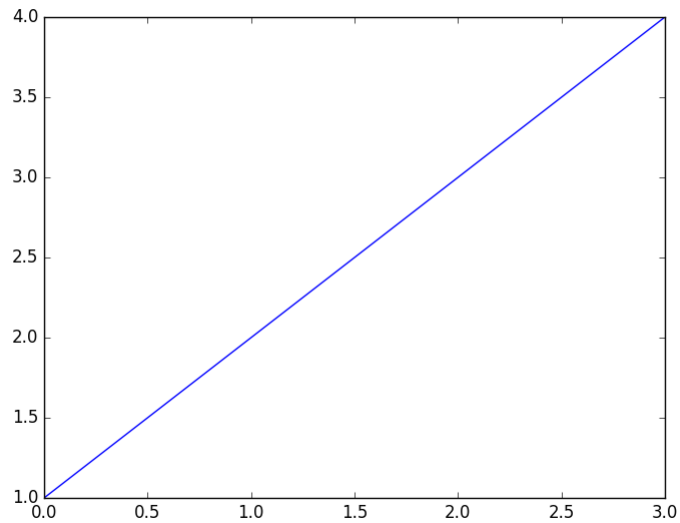
This is just a short introduction to the `matplotlib` plotting package. Its capabilities and customizations are described at length in the [project's webpage](#), the [Beginner's Guide](#), the `matplotlib.pyplot` [tutorial](#), and the `matplotlib.pyplot` [documentation](#). (Check in particular the [specific documentation](#) of `pyplot.plot`).

### H.2 Basic Usage – `pyplot.plot`

Simple use of `matplotlib` is straightforward:

```
>>> from matplotlib import pyplot as plt
>>> plt.plot([1,2,3,4])
[<matplotlib.lines.Line2D at 0x7faa8d9ba400>]
>>> plt.show()
```

If you run this code in the interactive Python interpreter, you should get a plot like this:



Two things to note from this plot:

- `pyplot.plot` assumed our single data list to be the *y*-values;
- in the absence of an *x*-values list, `[0, 1, 2, 3]` was used instead.

---

**Note:** `pyplot` is commonly used abbreviated as `plt`, just as `numpy` is commonly abbreviated as `np`. The remainder of this chapter uses the abbreviated form.

---

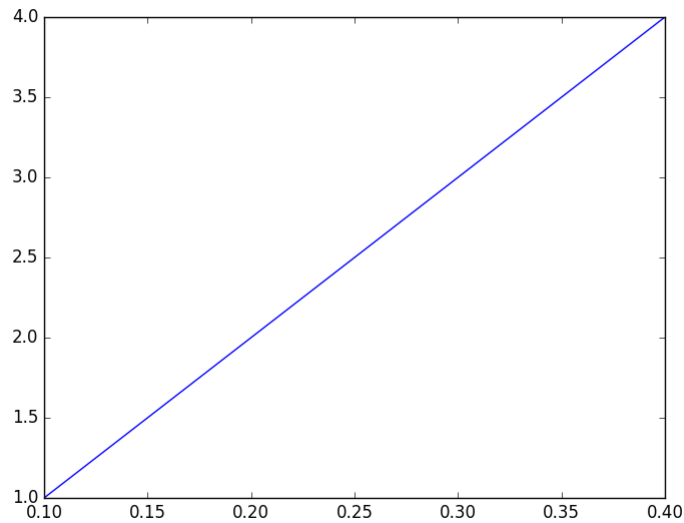
---

**Note:** Enhanced interactive python interpreters such as IPython can automate some of the plotting calls for you. For instance, you can run `%matplotlib` in IPython, after which you no longer need to run `plt.show` everytime when calling `plt.plot`. For simplicity, `plt.show` will also be left out of the remainder of these examples.

---

If you pass two lists to `plt.plot` you then explicitly set the *x* values:

```
>>> plt.plot([0.1, 0.2, 0.3, 0.4], [1, 2, 3, 4])
```

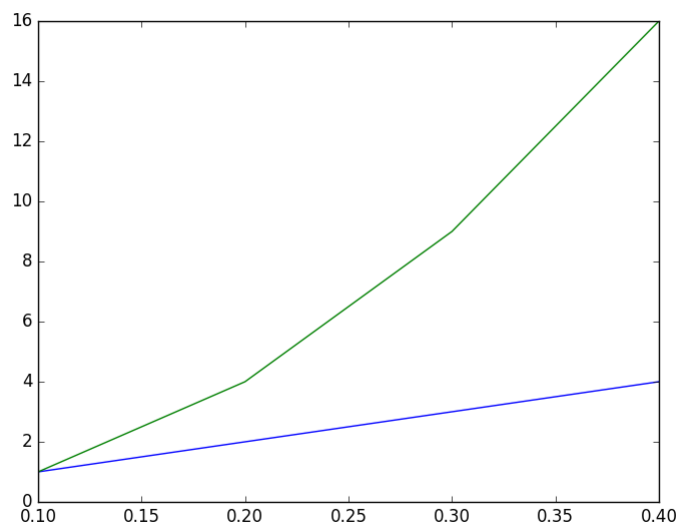


Understandably, if you provide two lists their lengths must match:

```
>>> plt.plot([0.1, 0.2, 0.3, 0.4], [1, 2, 3, 4, 5])
ValueError: x and y must have same first dimension
```

To plot multiple curves simply call `plt.plot` with as many *x-y* list pairs as needed:

```
>>> plt.plot([0.1, 0.2, 0.3, 0.4], [1, 2, 3, 4],
             [0.1, 0.2, 0.3, 0.4], [1, 4, 9, 16])
```

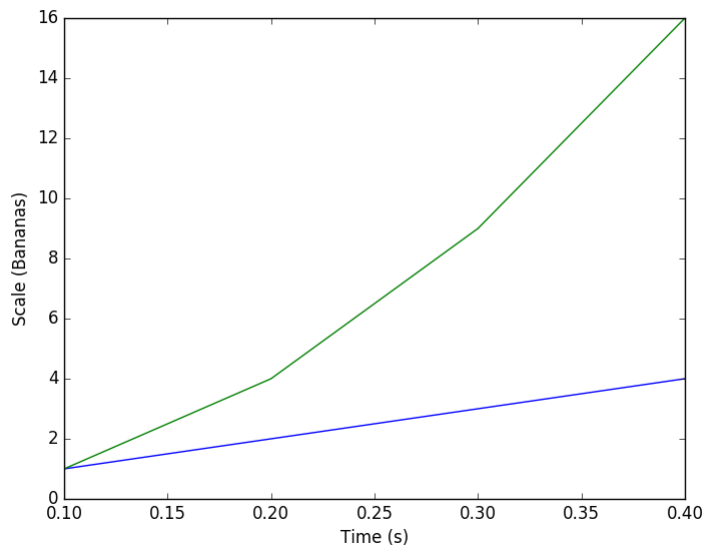


Alternatively, more plots may be added by repeatedly calling `plt.plot`. The following code snippet produces the same plot as the previous code example:

```
>>> plt.plot([0.1, 0.2, 0.3, 0.4], [1, 2, 3, 4])
>>> plt.plot([0.1, 0.2, 0.3, 0.4], [1, 4, 9, 16])
```

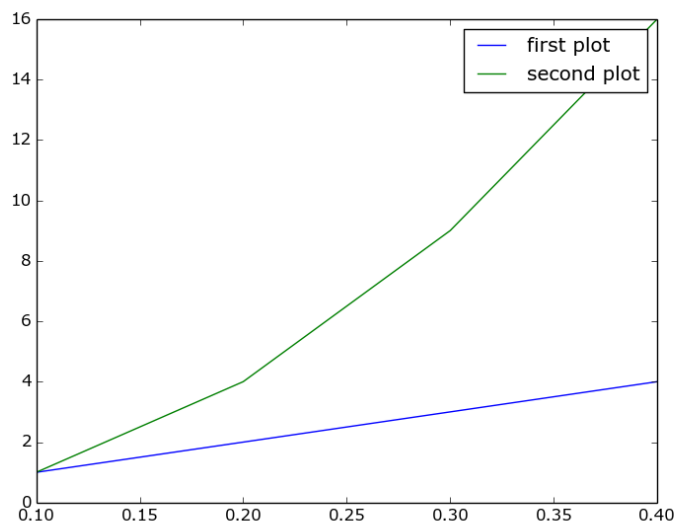
Adding information to the plot axes is straightforward to do:

```
>>> plt.plot([0.1, 0.2, 0.3, 0.4], [1, 2, 3, 4])
>>> plt.plot([0.1, 0.2, 0.3, 0.4], [1, 4, 9, 16])
>>> plt.xlabel("Time (s)")
>>> plt.ylabel("Scale (Bananas)")
```



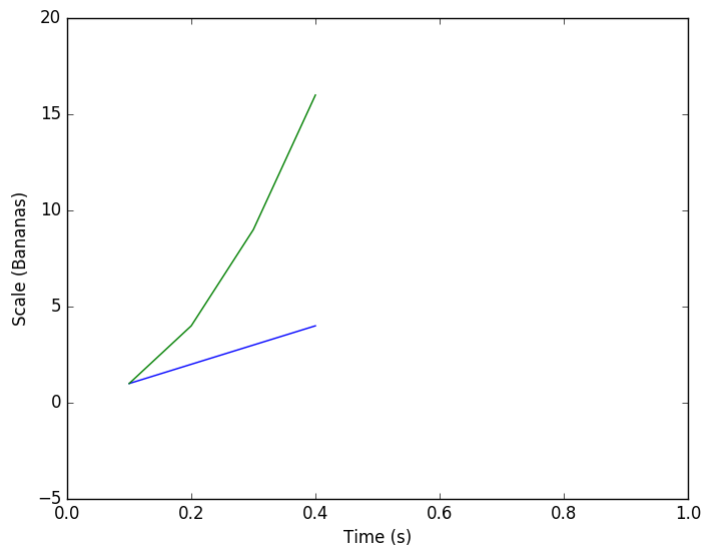
Also, adding an legend is rather simple:

```
>>> plt.plot([0.1, 0.2, 0.3, 0.4], [1, 2, 3, 4], label='first plot')
>>> plt.plot([0.1, 0.2, 0.3, 0.4], [1, 4, 9, 16], label='second plot')
>>> plt.legend()
```



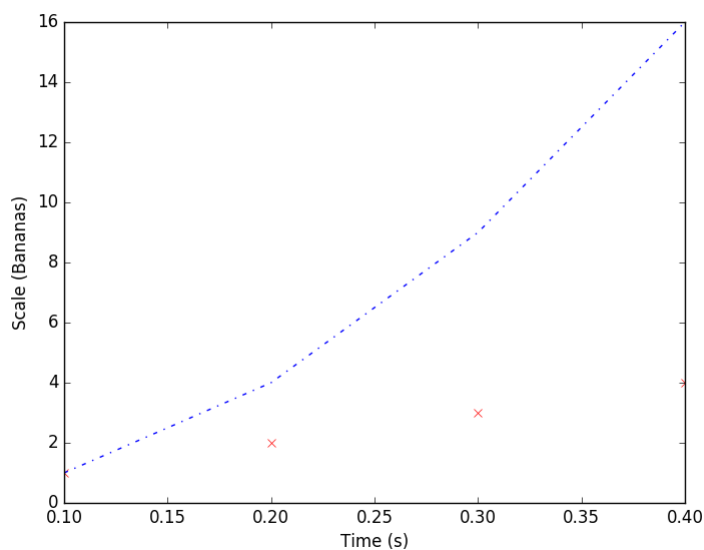
And adjusting axis ranges can be done by calling `plt.xlim` and `plt.ylim` with the lower and higher limits for the respective axes.

```
>>> plt.plot([0.1, 0.2, 0.3, 0.4], [1, 2, 3, 4])
>>> plt.plot([0.1, 0.2, 0.3, 0.4], [1, 4, 9, 16])
>>> plt.xlabel("Time (s)")
>>> plt.ylabel("Scale (Bananas)")
>>> plt.xlim(0, 1)
>>> plt.ylim(-5, 20)
```



In addition to *x* and *y* data lists, `plt.plot` can also take strings that define the plotting style:

```
>>> plt.plot([0.1, 0.2, 0.3, 0.4], [1, 2, 3, 4], 'rx')
>>> plt.plot([0.1, 0.2, 0.3, 0.4], [1, 4, 9, 16], 'b-.')
>>> plt.xlabel("Time (s)")
>>> plt.ylabel("Scale (Bananas)")
```



The style strings, one per *x*–*y* pair, specify color and shape: ‘*rx*’ stands for red crosses, and ‘*b-.*’ stands for blue dash-point line. Check the [documentation](#) of `pyplot.plot` for the list of colors and shapes.

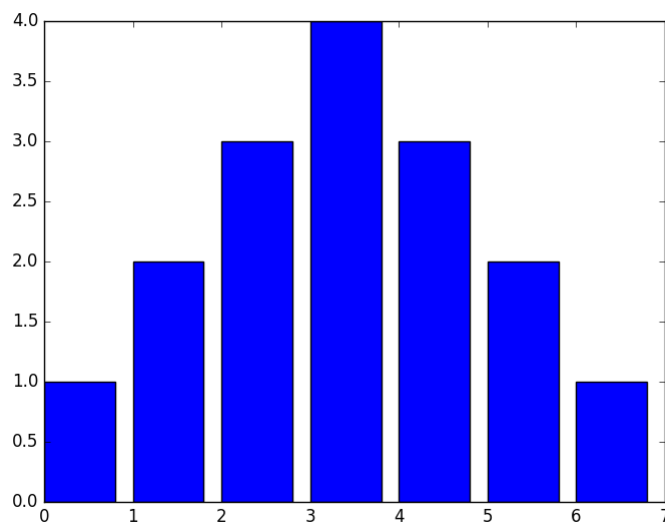
Finally, `plt.plot` can also, conveniently, take numpy arrays as its arguments.

## H.3 More plots

While `plt.plot` can satisfy basic plotting needs, `matplotlib` provides many more plotting functions. Below we try out the `plt.bar` function, for plotting bar charts. The full list of plotting functions can be found in the [matplotlib.pyplot documentation](#).

Bar charts can be plotted using `plt.bar`, in a similar fashion to `plt.plot`:

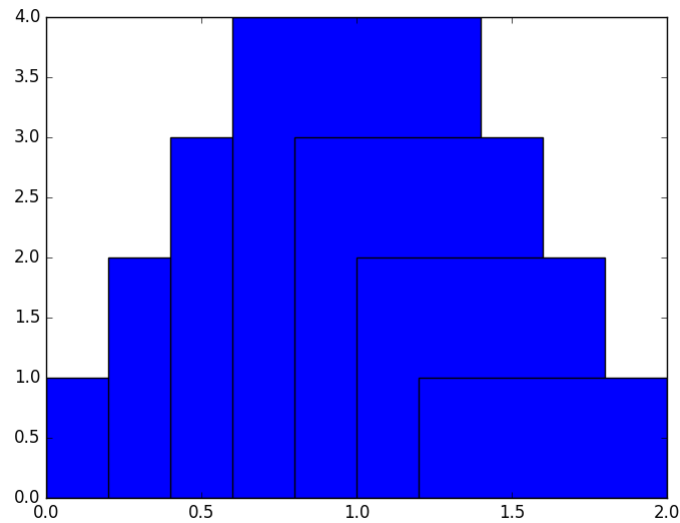
```
>>> plt.bar(range(7), [1, 2, 3, 4, 3, 2, 1])
```



Note, however, that contrary to `plt.plot` you must always specify *x* and *y* (which correspond, in bar chart terms to the left bin edges and the bar heights). Also note that you can only plot one chart per call. For multiple, overlapping charts you'll need to call `plt.bar` repeatedly.

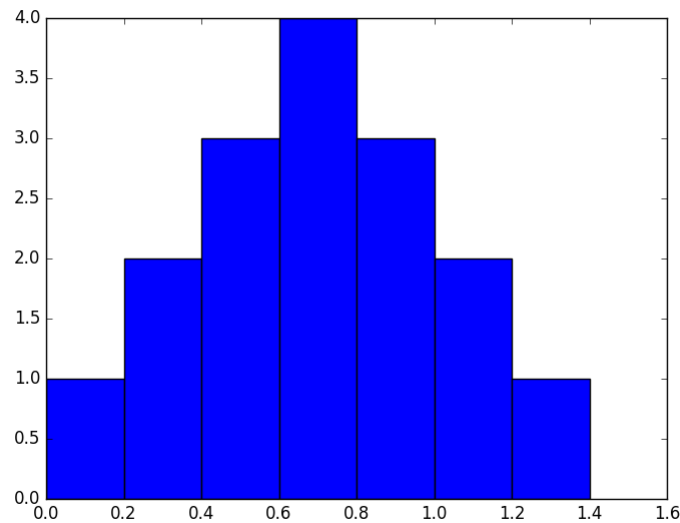
One of the optional arguments to `plt.bar` is `width`, which lets you specify the width of the bars. Its default of 0.8 might not be the most suited for all cases, especially when the *x* values are small:

```
>>> plt.bar(numpy.arange(0., 1.4, .2), [1, 2, 3, 4, 3, 2, 1])
```



Specifying narrower bars gives us a much better result:

```
>>> plt.bar(numpy.arange(0., 1.4, .2), [1, 2, 3, 4, 3, 2, 1], width=0.2)
```



Sometimes you will want to compare a function to your measured data; for example when you just fitted a function. Of course this is possible with matplotlib. Let's say we fitted a quadratic function to the first 10 prime numbers, and want to check how good our fit matches our data.

```
1 import matplotlib.pyplot as plt
2
3 def found_fit(x):
4     return 0.388 * x**2 # Found with sympfit.
5
6 x_data = list(range(10))
7 y_data = [2, 3, 5, 7, 11, 13, 17, 19, 23, 29]
```

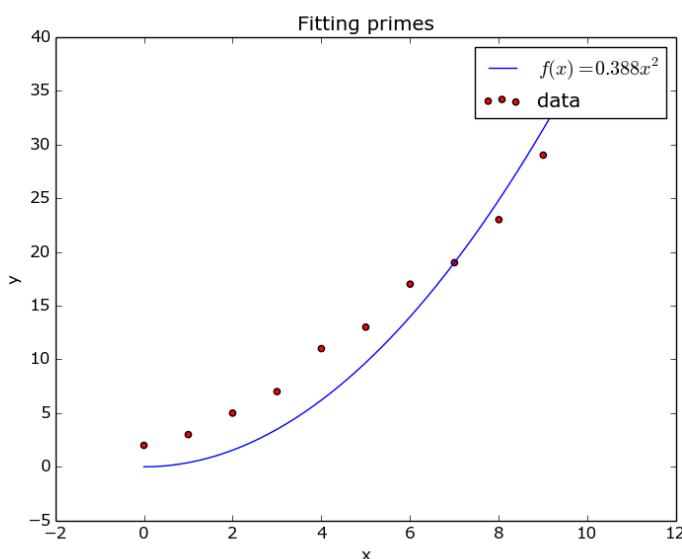
(continues on next page)



(continued from previous page)

```
8
9 x_func = np.linspace(0, 10, 50)
10 # numpy will do the right thing and evaluate found_fit for all elements
11 y_func = found_fit(x_func)
12
13 # From here the plotting starts
14
15 plt.scatter(x_data, y_data, c='r', label='data')
16 plt.plot(x_func, y_func, label='$f(x) = 0.388 x^2$')
17 plt.xlabel('x')
18 plt.ylabel('y')
19 plt.title('Fitting primes')
20 plt.legend()
21 plt.show()
```

We made the scatter plot red by passing it the keyword argument `c='r'`; `c` stands for colour, `r` for red. In addition, the label we gave to the `plot` statement is in *LaTeX* format, making it very pretty indeed. It's not a great fit, but that's besides the point here.

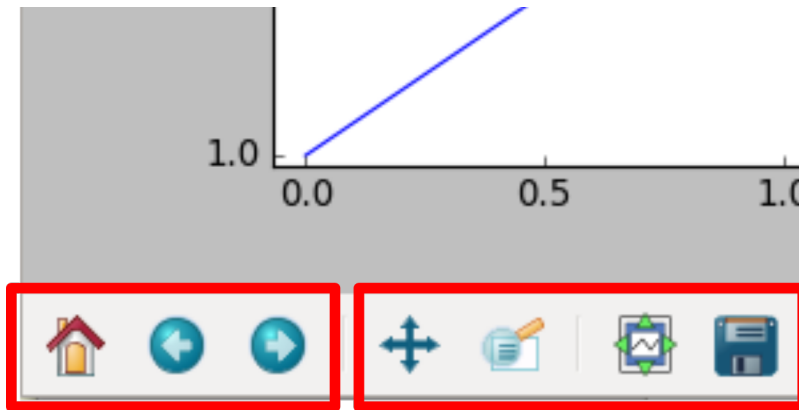


## H.4 Interactivity and saving to file

If you tried out the previous examples using a Python/IPython console you probably got for each plot an interactive window. Through the four rightmost buttons in this window you can do a number of actions:

- Pan around the plot area;
- Zoom in and out;
- Access interactive plot size control;
- Save to file.

The three leftmost buttons will allow you to navigate between different plot views, after zooming/panning.



As explained above, saving to file can be easily done from the interactive plot window. However, the need might arise to have your script write a plot directly as an image, and not bring up any interactive window. This is easily done by calling `plt.savefig`:

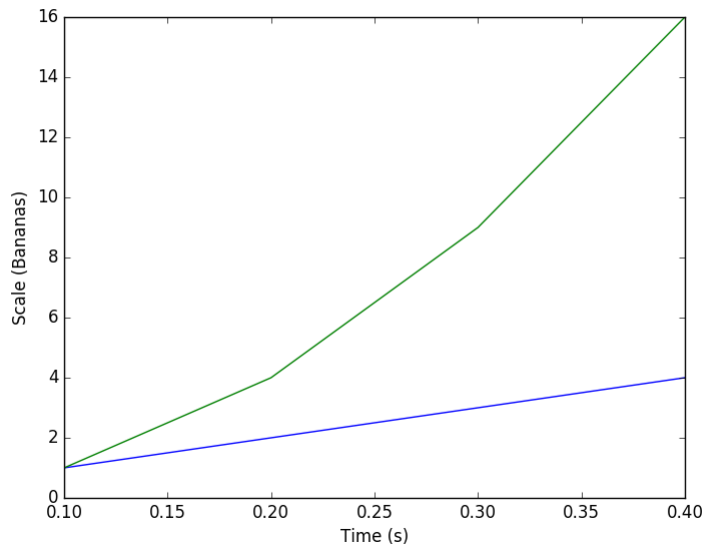
```
>>> plt.plot([0.1, 0.2, 0.3, 0.4], [1, 2, 3, 4], 'rx')
>>> plt.plot([0.1, 0.2, 0.3, 0.4], [1, 4, 9, 16], 'b-.')
>>> plt.xlabel("Time (s)")
>>> plt.ylabel("Scale (Bananas)")
>>> plt.savefig('the_best_plot.pdf')
```

**Note:** When saving a plot, you'll want to choose a **vector format** (either pdf, ps, eps, or svg). These are resolution-independent formats and will yield the best quality, even if printed at very large sizes. Saving as png should be avoided, and saving as jpg should be avoided even more.

## H.5 Multiple figures

With this groundwork out of the way, we can move on to some more advanced matplotlib use. It is also possible to use it in an object-oriented manner, which allows for more separation between several plots and figures. Let's say we have two sets of data we want to plot next to each other, rather than in the same figure. Matplotlib has several layers of organisation: first, there's an `Figure` object, which basically is the window your plot is drawn in. On top of that, there are `Axes` objects, which are your separate graphs. It is perfectly possible to have multiple (or no) `Axes` in one `Figure`. We'll explain the `add_subplot` method a bit later. For now, it just creates an `Axis` instance.

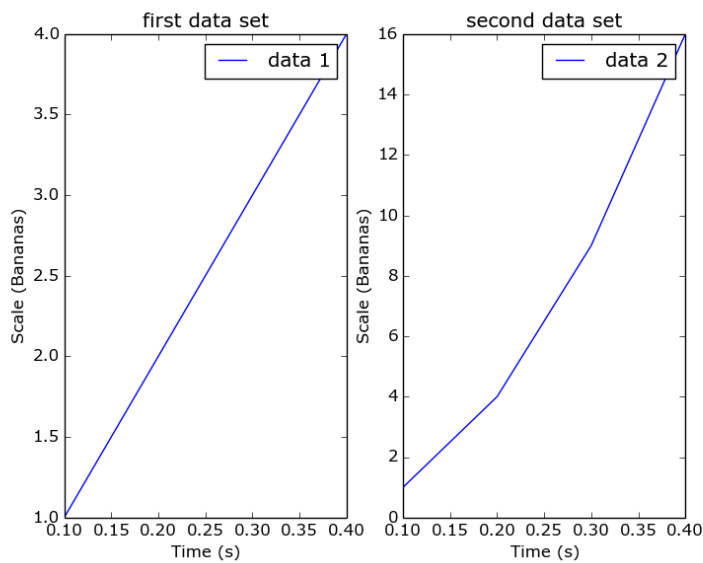
```
1 import matplotlib.pyplot as plt
2
3 x_data = [0.1, 0.2, 0.3, 0.4]
4 y_data = [1, 2, 3, 4]
5
6 fig = plt.figure()
7 ax = fig.add_subplot(1, 1, 1)
8 ax.plot([0.1, 0.2, 0.3, 0.4], [1, 2, 3, 4])
9 ax.plot([0.1, 0.2, 0.3, 0.4], [1, 4, 9, 16])
10 ax.set_xlabel('Time (s)')
11 ax.set_ylabel('Scale (Bananas)')
12
13 plt.show()
```



This example also neatly highlights one of Matplotlib's shortcomings: the API is highly inconsistent. Where we could do `xlabel()` before, we now need to do `set_xlabel()`. In addition, we can't show the figures one by one (i.e. `fig.show()`); instead we can only show them all at the same time with `plt.show()`.

Now, we want to make multiple plots next to each other. We do that by calling `plot` on two different axes:

```
1  x_data1 = [0.1, 0.2, 0.3, 0.4]
2  y_data1 = [1, 2, 3, 4]
3
4  x_data2 = [0.1, 0.2, 0.3, 0.4]
5  y_data2 = [1, 4, 9, 16]
6
7  fig = plt.figure()
8  ax1 = fig.add_subplot(1, 2, 1)
9  ax2 = fig.add_subplot(1, 2, 2)
10 ax1.plot(x_data1, y_data1, label='data 1')
11 ax2.plot(x_data2, y_data2, label='data 2')
12 ax1.set_xlabel('Time (s)')
13 ax1.set_ylabel('Scale (Bananas)')
14 ax1.set_title('first data set')
15 ax1.legend()
16 ax2.set_xlabel('Time (s)')
17 ax2.set_ylabel('Scale (Bananas)')
18 ax2.set_title('second data set')
19 ax2.legend()
20
21 plt.show()
```



The `add_subplot` method returns an `Axis` instance and takes three arguments: the first is the number of rows to create; the second is the number of columns; and the last is which plot number we add right now. So in common usage you will need to call `add_subplot` once for every axis you want to make with the same first two arguments. What would happen if you first ask for one row and two columns, and for two rows and one column in the next call?

## H.6 Exercises

1. Plot a dashed line.
2. Search the matplotlib documentation, and plot a line with plotmarkers on all it's datapoints. You can do this with just one call to `plt.plot`.