

Dive Into Python

Chapter 1. Installing

- 1.1 Which is right for you?
- 1.2 Python on Windows
- 1.3 Python on Mac OS X
- 1.4 Python on Mac OS 9
- 1.5 Python on RedHat Linux
- 1.6 Python on Debian GNU/Linux
- 1.7 Python Installation from Source
- 1.8 Interactive Shell
- 1.9. Summary

Chapter 2. Your Program

- 2.1 Diving
- 2.2 Fundamentals
- 2.3 Functions
- 2.4.1 Everything
- 2.5 Indenting
- 2.6 Modules

Chapter 3. Native

- 3.1 Diving
- 3.2 Introducing
- 3.3 Tuples
- 3.4 Variables
- 3.5 String Formatting
- 3.6 Mapping
- 3.7 Joining and Splitting Strings
- 3.8. Summary

Chapter 4. Introspection

- 4.1 Diving
- 4.2 Optional and Named Arguments
- 4.3 dir, isinstance, and Other Built-In Functions
- 4.4 Object References With getattr
- 4.5 Filtering
- 4.6 The Nature of and and or
- 4.7 Built-In Functions
- 4.8 Putting Together
- 4.9. Summary

Chapter 5. Object Orientation

- 5.1 Diving
- 5.2 Importing from module import
- 5.3 Class Defining
- 5.4 Class Instantiating
- 5.5 Example: A Wrapper Class
- 5.6 Special Methods
- 5.7 Special Class Methods
- 5.8 Class Attributes

Table of Contents

Chapter 5: Orientation

- 5.1. Private
- 5.10. Summary

Chapter 6: File Handling

- 6.1. Opening
- 6.2. Working with File Objects
- 6.3. File Objects
- 6.4. Closing
- 6.5. Working with Directories
- 6.6. Putting it together
- 6.7. Summary

Chapter 7: Regular Expressions

- 7.1. Diving
- 7.2. Using Street Addresses
- 7.3. Using Roman Numerals
- 7.4. Using Syntax
- 7.5. Using Expressions
- 7.6. Parsing Phone Numbers
- 7.7. Summary

Chapter 8: HTML

- 8.1. Diving
- 8.2. Multiplying
- 8.3. From HTML documents
- 8.4. HTML processor.py
- 8.5. locals and globals
- 8.6. Dictionary-based
- 8.7. Quoting
- 8.8. Multiplying
- 8.9. Putting it together
- 8.10. Summary

Chapter 9: XML

- 9.1. Diving
- 9.2. Packages
- 9.3. Parsing
- 9.4. Unicode
- 9.5. Searching
- 9.6. Accessing attributes
- 9.7. Segue

Chapter 10: Scripts

- 10.1. Abstracting
- 10.2. Standard and error
- 10.3. Clackings
- 10.4. Finding a node
- 10.5. Creating handlers by node type
- 10.6. Handling arguments
- 10.7. Putting it together

Table of Contents

Chapter 10. Scripts

10.8. Summary

Chapter 11. HTTP

11.1. Diving
11.2. How to fetch data over HTTP
11.3. HTTP features
11.4. HTTP debugging services
11.5. User-Agent
11.6. Modifying and ETag
11.7. Handling
11.8. Handling data
11.9. Putting it all together
11.10. Summary

Chapter 12. SOAP

12.1. Diving
12.2. SOAP Libraries
12.3. First with SOAP
12.4. Web Services
12.5. Introducing
12.6. Web Services with WSDL
12.7. Searching
12.8. SOAP Web Services
12.9. Summary

Chapter 13. Unit

13.1. Introduction
13.2. Diving
13.3. Making
13.4. Testing
13.5. Testing
13.6. Testing

Chapter 14. Test-First

14.1. man.py, stage 1
14.2. man.py, stage 2
14.3. man.py, stage 3
14.4. man.py, stage 4
14.5. man.py, stage 5

Chapter 15. Refactoring

15.1. Handling
15.2. Handling requirements
15.3. Refactoring
15.4. Postscript
15.5. Summary

Chapter 16. Functional

16.1. Diving
16.2. Finding

Table of Contents

Chapter 16. Functional

- 16.1. Lists
- 16.2. Mapping
- 16.3. Data structures
- 16.4. Dynamic objects
- 16.5. All together
- 16.8. Summary

Chapter 17. Dynamic

- 17.1. Diving
- 17.2. al.py, stage 1
- 17.3. al.py, stage 2
- 17.4. al.py, stage 3
- 17.5. al.py, stage 4
- 17.6. al.py, stage 5
- 17.7. al.py, stage 6
- 17.8. Summary

Chapter 18. Performance

- 18.1. Diving
- 18.2. Using Module
- 18.3. Optimizations
- 18.4. Optimizing loops
- 18.5. Optimizing
- 18.6. Manipulation
- 18.7. Summary

Appendix A. Further

Appendix B. Review

Appendix C. Tips

Appendix D. List

Appendix E. Revision

Appendix F. About

Appendix G. GNU License

- G.0. Preamble
- G.1. Applicability
- G.2. Copying
- G.3. Copying
- G.4. Modifications
- G.5. Combining
- G.6. Collections
- G.7. Independent works
- G.8. Translation
- G.9. Termination
- G.10. Future of this license

Table of Contents

Appendix A. Creative Commons Attribution License

How to use this License for your documents

Appendix H. Python

How to use Python

How to use Python for accessing or otherwise using Python

Dive Into Python

20 May 2004

Copyright © 2000, 2001, 2002, 2003, 2004 Mark Pilgrim (mailto:mark@diveintopython.org)

This book lives at <http://diveintopython.org/>. If you're reading it somewhere else, you may not have the latest version.

Permission is granted to copy, distribute, and/or modify this document under the terms of the GNU Free Documentation License, Version 1.1.

The example programs in this book are free software; you can redistribute and/or modify them under the terms of the Python license as published by the Free Software Foundation.

Chapter 1. Installing Python

Welcome to Python. Let's dive in. In this chapter, you'll install the version of Python that's right for you.

1.1. Which Python is right for you?

The first thing you need to do with Python is install it. Or do you?

If you're using an account on a hosted server, your ISP may have already installed Python. Most popular Linux distributions come with Python.

Windows does not come with any version of Python, but don't despair! There are several ways to point-and-click your way to Python on Windows.

As you can see already, Python runs on a great many operating systems. The full list includes Windows, Mac OS, Mac OS X, and all varieties of Unix.

What's more, Python programs written on one platform can, with a little care, run on *any* supported platform. For instance, I regularly develop Python programs on Linux.

So back to the question that started this section, "Which Python is right for you?" The answer is whichever one runs on the computer you already have.

1.2. Python on Windows

On Windows, you have a couple choices for installing Python.


ActiveState makes a Windows installer for Python called ActivePython, which includes a complete version of Python, an IDE with a Python interpreter, and a lot of other goodies.

ActivePython is freely downloadable, although it is not open source. It is the IDE I used to learn Python, and I recommend you try it unless you have a good reason not to.

The second option is the "official" Python installer, distributed by the people who develop Python itself. It is freely downloadable and open source.

Procedure 1.1. Option 1: Installing ActivePython




Here is the procedure for installing ActivePython:

1.
Download ActivePython from <http://www.activestate.com/Products/ActivePython/>.
2.
If you are using Windows 95, Windows 98, or Windows ME, you will also need to download and install Windows Installer 2.0 (<http://www.microsoft.com/windows/win98/active/activepython.htm>).
3.
Double-click the installer, .
4.
Step through the installer program.
5.
If space is tight, you can do a custom installation and deselect the documentation, but I don't recommend this unless you absolutely have to.
6.
After the installation is complete, close the installer and choose Start->Programs->ActiveState ActivePython 2.2->PythonWin IDE.



>

Procedure 1.2. Option 2: Installing Python from Python.org (<http://www.python.org/>)

1.
Download the latest Python Windows installer by going to <http://www.python.org/ftp/python/> and selecting the highest version number.
2.
Double-click the installer, The name will depend on the version of Python available when you read this.
3.
Step through the installer program.
4.
If disk space is tight, you can deselect the HTMLHelp file, the utility scripts () and/or the test suite ().
5.
If you do not have administrative rights on your machine, you can select Advanced Options, then choose Non-Admin Install. This just works.
6.
After the installation is complete, close the installer and select Start->Programs->Python 2.3->IDLE (Python GUI). You'll see some



1.3. Python on Mac OS X



On Mac OS X, you have two choices for installing Python: install it, or don't install it. You probably want to install it.


Mac OS X 10.2 and later comes with a command-line version of Python preinstalled. If you are comfortable with the command line, you can


Rather than using the preinstalled version, you'll probably want to install the latest version, which also comes with a graphical interactive shell.

Procedure 1.3. Running the Preinstalled Version of Python on Mac OS X

To use the preinstalled version of Python, follow these steps:

1.
Open the  folder.
2.
Open the  folder.

Double-click  to open a terminal window and get to a command line.³







4.
Type  at the command prompt.

Try it out:



Procedure 1.4. Installing the Latest Version of Python on Mac OS X

Follow these steps to download and install the latest version of Python:

1.
Download the  disk image from <http://homepages.cwi.nl/~jack/macpython/download.html>.
2.
If your browser has not already done so, double-click  to mount the disk image on your desktop.
3.
Double-click the installer, .
4.
The installer will prompt you for your administrative username and password.
5.
Step through the installer program.
6.
After installation is complete, close the installer and open the  folder.
7.
Open the  folder
8.
Double-click  to launch Python.

The MacPython IDE should display a splash screen, then take you to the interactive shell. If the interactive shell does not appear, select Window



Note that once you install the latest version, the pre-installed version is still present. If you are running scripts from the command line, you need to specify the version.


Example 1.1. Two versions of Python

```
python --help
python --help
python --help
python --help
python --help
python --help
python --help
python --help
python --help
python --help
```

1.4. Python on Mac OS 9

Mac OS 9 does not come with any version of Python, but installation is very simple, and there is only one choice.

Follow these steps to install Python on Mac OS 9:

1.
Download the file from <http://homepages.cwi.nl/~jack/macpython/download.html>.
2.
If your browser does not decompress the file automatically, double-click the file to decompress the file with Stuffit Expander.
3.
Double-click the installer, .
4.
Step through the installer program.
5.
After installation is complete, close the installer and open the folder.
6.
Open the folder.
7.
Double-click the icon to launch Python.

The MacPython IDE should display a splash screen, and then take you to the interactive shell. If the interactive shell does not appear, select View > Shell.

```
python --help
python --help
python --help
python --help
python --help
python --help
python --help
python --help
python --help
python --help
```




1.7. Python Installation from Source

If you prefer to build from source, you can download the Python source code from <http://www.python.org/ftp/python/>. Select the highest version available.

Example 1.4. Installing from source





1.8. The Interactive Shell

Now that you have Python installed, what's this interactive shell thing you're running?

It's like this: Python leads a double life. It's an interpreter for scripts that you can run from the command line or run like applications, by dou

Launch the Python interactive shell in whatever way works on your platform, and let's dive in with the steps shown here:

Example 1.5. First Steps in the Interactive Shell



- 1 The Python interactive shell can evaluate arbitrary Python expressions, including any basic arithmetic expression.
- 2 The interactive shell can execute arbitrary Python statements, including the **print** statement.
- 3 You can also assign values to variables, and the values will be remembered as long as the shell is open (but not any longer than that).

1.9. Summary

You should now have a version of Python installed that works for you.

Depending on your platform, you may have more than one version of Python installed. If so, you need to be aware of your paths. If simply ty

Congratulations, and welcome to Python.

Chapter 2. Your First Python Program

You know how other books go on and on about programming fundamentals and finally work up to building a complete, working program? L

2.1. Diving in

Here is a complete, working Python program.

It probably makes absolutely no sense to you. Don't worry about that, because you're going to dissect it line by line. But read through it first

Example 2.1.

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5>).



```
  
  
  
  
  
  
  
  
  
  
  
  
  

```

Now run this program and see what happens.

In the ActivePython IDE on Windows,  you can run the Python program you're editing by choosing File->Run... (**Ctrl-R**). Output is displayed

In the Python IDE on Mac OS,  you can run a Python program with Python->Run window... (**Cmd-R**), but there is an important option you m

On UNIX-compatible systems  (including Mac OS X), you can run a Python program from the command line:

The output of  will look like this:

```

```

2.2. Declaring Functions

Python has functions like most other languages, but it does not have separate header files like C++ or  sections like Pascal. When you need a

```

```

Note that the keyword `def` starts the function declaration, followed by the function name, followed by the arguments in parentheses. Multiple arguments are separated by commas.

Also note that the function doesn't define a return datatype. Python functions do not specify the datatype of their return value; they don't even have a return statement.

In Visual Basic, functions (that return a value) start with `Function` and subroutines (that do not return a value) start with `Sub`. There are no subroutines in Python.

The argument, `name`, doesn't specify a datatype. In Python, variables are never explicitly typed. Python figures out what type a variable is and keeps track of it.

In Java, C++, and other statically-typed languages, you must specify the datatype of the function return value and each function argument. In Python, you don't.

2.2.1. How Python's Datatypes Compare to Other Programming Languages

An erudite reader sent me this explanation of how Python compares to other programming languages:

statically typed language

A language in which types are fixed at compile time. Most statically typed languages enforce this by requiring you to declare variables and their types.

dynamically typed language

A language in which types are discovered at execution time; the opposite of statically typed. VBScript and Python are dynamically typed.

strongly typed language

A language in which types are always enforced. Java and Python are strongly typed. If you have an integer, you can't treat it as a string.

weakly typed language

A language in which types may be ignored; the opposite of strongly typed. VBScript is weakly typed. In VBScript, you can treat an integer as a string.

So Python is both *dynamically typed* (because it doesn't use explicit datatype declarations) and *strongly typed* (because once a variable has a type, it can't be treated as a different type).

2.3. Documenting Functions

You can document a Python function by giving it a docstring.

Example 2.2. Defining the `len` Function's Docstring

```
def len(s):  
    """  
    Return the length of the string s.  
    """  
    return len(s)
```

Triple quotes signify a multi-line string. Everything between the start and end quotes is part of a single string, including carriage returns and newlines.

Triple quotes are also an easy way to define a string with both single and double quotes, like `"' in Perl`.

Everything between the triple quotes is the function's docstring, which documents what the function does. A docstring, if it exists, must be the first thing defined in the function.

Many Python IDEs use the docstring to provide context-sensitive documentation, so that when you type a function name, its docstring appears as a tooltip. This is a handy feature.

Further Reading on Documenting Functions

PEP 257 (<http://www.python.org/peps/pep-0257.html>) defines **conventions**.

Python Style Guide (<http://www.python.org/doc/essays/styleguide.html>) discusses how to write a good **code**.

Python Tutorial (<http://www.python.org/doc/current/tut/tut.html>) discusses conventions for spacing in **code** (<http://www.python.org/doc/current/tut/tut.html>).

2.4. Everything Is an Object

In case you missed it, I just said that Python functions have attributes, and that those attributes are available at runtime.

A function, like everything else in Python, is an object.

Open your favorite Python IDE and follow along:

Example 2.3. Accessing the **Function's** **Attributes**

```
1 import sys
2
3 def main():
4     print(sys.argv[1])
5
6 if __name__ == '__main__':
7     main()
```

- 1 The first line imports the **sys** module -- a chunk of code that you can use interactively, or from a larger Python program. (You can also import modules from other Python programs.)
- 2 When you want to use functions defined in imported modules, you need to include the module name. So you can't just say **main()**; it must be **main()**.
- 3 Instead of calling the function as you would expect to, you asked for one of the function's attributes, **main.__doc__**.

In Python is like in Perl. Once you **import** a Python module, you access its functions with **module.function()**. Once you **require** a Perl module, you access its functions with **module->function()**.

2.4.1. The Import Search Path

Before you go any further, I want to briefly mention the library search path. Python looks in several places when you try to import a module.

Example 2.4. Import Search Path

```
1 import sys
2
3 def main():
4     print(sys.argv[1])
5
6 if __name__ == '__main__':
7     main()
```

- 1 Importing the **sys** module makes all of its functions and attributes available.
- 2 **sys.path** is a list of directory names that constitute the current search path. (Yours will look different, depending on your operating system, where Python is installed.)
- 3 Actually, I lied; the truth is more complicated than that, because not all modules are stored as **files**. Some, like the **sys** module, are "built-in" modules.
- 4 You can add a new directory to Python's search path at runtime by appending the directory name to **sys.path** and then Python will look in that directory for modules.

2.4.2. What's an Object?

Everything in Python is an object, and almost everything has attributes and methods. All functions have a built-in attribute `__doc__` which returns the

Still, this begs the question. What is an object? Different programming languages define "object" in different ways. In some, it means that *all*

This is so important that I'm going to repeat it in case you missed it the first few times: *everything in Python is an object*. Strings are objects.

Further Reading on Objects

Python Reference Manual (<http://www.python.org/doc/current/ref/>) explains exactly what it means to say that everything in Python is an object. *eff-bot* (<http://www.effbot.org/guides/>) summarizes Python objects (<http://www.effbot.org/guides/python-objects.htm>).

2.5. Indenting Code

Python functions have no explicit `begin` and no curly braces to mark where the function code starts and stops. The only delimiter is a colon (`:`) at the

Example 2.5. Indenting the `def` Function

```
def fib(n):  
    """Fibonacci sequence: 0 1 1 2 3 5 8 13 21 34 55 89 144 ...  
    the n<sup>th</sup> Fibonacci number is the sum of the two preceding numbers.  
    """  
    if n < 0:  
        raise ValueError("Fibonacci sequence: n must be a non-negative integer")  
    if n < 2:  
        return n  
    a, b = 0, 1  
    for i in range(n - 1):  
        a, b = b, a + b  
    return b
```

Code blocks are defined by their indentation. By "code block", I mean functions, `if` statements, `for` loops, `while` loops, and so forth. Indenting starts a block.

Example 2.6, "if Statements" shows an example of code indentation with `if` statements.

Example 2.6. `if` Statements

```
def fib(n):  
    """Fibonacci sequence: 0 1 1 2 3 5 8 13 21 34 55 89 144 ...  
    the n<sup>th</sup> Fibonacci number is the sum of the two preceding numbers.  
    """  
    if n < 0:  
        raise ValueError("Fibonacci sequence: n must be a non-negative integer")  
    if n < 2:  
        return n  
    a, b = 0, 1  
    for i in range(n - 1):  
        a, b = b, a + b  
    return b
```

- 1 This is a function named `fib` that takes one argument, `n`. All the code within the function is indented.
- 2 Printing to the screen is very easy in Python, just use `print`. Statements can take any data type, including strings, integers, and other native Python objects.
- 3 `if` statements are a type of code block. If the `if` expression evaluates to true, the indented block is executed, otherwise it falls to the `else` block.
- 4 Of course `if` and `else` blocks can contain multiple lines, as long as they are all indented the same amount. This `if` block has two lines of code indented.

After some initial protests and several snide analogies to Fortran, you will make peace with this and start seeing its benefits. One major benefit is that it makes the code more readable.

Python uses carriage returns to separate statements and a colon and indentation to separate code blocks. C++ and Java use semicolons to separate statements.

Further Reading on Code Indentation

Python Reference Manual (<http://www.python.org/doc/current/ref/>) discusses cross-platform indentation issues and shows various indentation styles. *Python Style Guide* (<http://www.python.org/doc/essays/styleguide.html>) discusses good indentation style.

2.6. Testing Modules

Python modules are objects and have several useful attributes. You can use this to easily test your modules as you write them. Here's an example:

```
def test():
```

Some quick observations before you get to the good stuff. First, parentheses are not required around the `if __name__ == '__main__':` expression. Second, the `if __name__ == '__main__':` statement is

Like C, Python uses `==` for comparison and `=` for assignment. Unlike C, Python does not support in-line assignment, so there's no chance of accidentally

So why is this particular `if __name__ == '__main__':` statement a trick? Modules are objects, and all modules have a built-in attribute `__name__`. A module's `__name__` depends on how you

```
    pass
```

```
if __name__ == '__main__':
```

```
    test()
```

Knowing this, you can design a test suite for your module within the module itself by putting it in this `if __name__ == '__main__':` statement. When you run the module

On MacPython, there is an additional step to make the `if __name__ == '__main__':` trick work. Pop up the module's options menu by clicking the black triangle in the upper left

Further Reading on Importing Modules

Python Reference Manual (<http://www.python.org/doc/current/ref/>) discusses the low-level details of importing modules (<http://www.python.org/doc/current/ref/>)

Chapter 3. Native Datatypes

You'll get back to your first Python program in just a minute. But first, a short digression is in order, because you need to know about dictionaries.

3.1. Introducing Dictionaries

One of Python's built-in datatypes is the dictionary, which defines one-to-one relationships between keys and values.

A dictionary in Python is like a hash in Perl. In Perl, variables that store hashes always start with a % character. In Python, variables can be named anything.

A dictionary in Python is like an instance of the Map class in Java.

A dictionary in Python is like an instance of the Dictionary object in Visual Basic.

3.1.1. Defining Dictionaries

Example 3.1. Defining a Dictionary

```
1 d = {'cat': 'feline', 'dog': 'canine'}
2
3 print(d['cat'])
4
5 print(d['dog'])
6
7 print(d.get('cat'))
8
9 print(d.get('cat', 'unknown'))
10
11 print(d.get('bird', 'unknown'))
```

- 1 First, you create a new dictionary with two elements and assign it to the variable `d`. Each element is a key-value pair, and the whole set is enclosed in curly braces.
- 2 `'cat'` is a key, and its associated value, referenced by `'feline'`, is the value.
- 3 `'dog'` is a key, and its associated value, referenced by `'canine'`, is the value.
- 4 You can get values by key, but you can't get keys by value. So `d['cat']` works, but `d['feline']` raises an exception, because `'feline'` is not a key.

3.1.2. Modifying Dictionaries

Example 3.2. Modifying a Dictionary

```
1 d = {'cat': 'feline', 'dog': 'canine'}
2
3 d['cat'] = 'kitten'
4
5 d['dog'] = 'puppy'
6
7 print(d)
```

- ❶ You can not have duplicate keys in a dictionary. Assigning a value to an existing key will wipe out the old value.
- ❷ You can add new key-value pairs at any time. This syntax is identical to modifying existing values. (Yes, this will annoy you someday.)

Note that the new element (key `1` value `1`) appears to be in the middle. In fact, it was just a coincidence that the elements appeared to be in order.

Dictionaries have no concept of order among elements. It is incorrect to say that the elements are "out of order"; they are simply unordered.

When working with dictionaries, you need to be aware that dictionary keys are case-sensitive.

Example 3.3. Dictionary Keys Are Case-Sensitive

```

1 d = {}
2 d['A'] = 1
3 d['a'] = 2
4 print(d)
5
6 d['A'] = 3
7 print(d)
8
9 d['a'] = 4
10 print(d)

```

- ❶ Assigning a value to an existing dictionary key simply replaces the old value with a new one.
- ❷ This is not assigning a value to an existing dictionary key, because strings in Python are case-sensitive, so `'A'` is not the same as `'a'`. This creates a new key-value pair.

Example 3.4. Mixing Datatypes in a Dictionary

```

1 d = {}
2 d[1] = 'one'
3 d[2] = 'two'
4 d['three'] = 3
5 print(d)
6
7 d[4] = {'five': 5}
8 print(d)
9
10 d[5] = {'six': 6}
11 print(d)

```

- ❶ Dictionaries aren't just for strings. Dictionary values can be any datatype, including strings, integers, objects, or even other dictionaries.
- ❷ Dictionary keys are more restricted, but they can be strings, integers, and a few other types. You can also mix and match key datatypes.

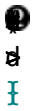
3.1.3. Deleting Items From Dictionaries

Example 3.5. Deleting Items from a Dictionary

```

1 d = {}
2 d[1] = 'one'
3 d[2] = 'two'
4 d[3] = 'three'
5 print(d)
6
7 del d[1]
8 print(d)
9
10 d.pop(2)
11 print(d)

```



- 1 Lets you delete individual items from a dictionary by key.
- 2 Deletes all items from a dictionary. Note that the set of empty curly braces signifies a dictionary without any items.

Further Reading on Dictionaries

How to Think Like a Computer Scientist (<http://www.ibiblio.org/obp/thinkCSpy/>) teaches about dictionaries and shows how to use them.
Python Knowledge Base (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>) has a lot of example code using dictionaries.
Python Cookbook (<http://www.activestate.com/ASPN/Python/Cookbook/>) discusses how to sort the values of a dictionary by key (http://www.python.org/doc/current/lib/).
Python Library Reference (<http://www.python.org/doc/current/lib/>) summarizes all the dictionary methods (<http://www.python.org/doc/current/lib/>).

3.2. Introducing Lists

Lists are Python's workhorse datatype. If your only experience with lists is arrays in Visual Basic or (God forbid) the datastore in Powerbuilder, you're in for a surprise.

A list in Python is like an array in Perl. In Perl, variables that store arrays always start with the @ character; in Python, variables can be named anything you like.

A list in Python is much more than an array in Java (although it can be used as one if that's really all you want out of life). A better analogy would be a vector in C++.

3.2.1. Defining Lists

Example 3.6. Defining a List



- 1 First, you define a list of five elements. Note that they retain their original order. This is not an accident. A list is an ordered set of elements.
- 2 A list can be used like a zero-based array. The first element of any non-empty list is always 0.
- 3 The last element of this five-element list is 4 because lists are always zero-based.

Example 3.7. Negative List Indices



- 1 A negative index accesses elements from the end of the list counting backwards. The last element of any non-empty list is always -1.
- 2 If the negative index is confusing to you, think of it this way: -1. So in this list, -1 is the last element.

Example 3.8. Slicing a List

- ```
1 # Create a list of five elements
2 my_list = ['a', 'b', 'c', 'd', 'e']
3
4 # Slice the list to get the first three elements
5 first_three = my_list[0:3]
6
7 # Print the original list and the slice
8 print(my_list)
9 print(first_three)
```
- 1 You can get a subset of a list, called a "slice", by specifying two indices. The return value is a new list containing all the elements of the original list between the two indices (including the element at the start index but not including the element at the end index).
  - 2 Slicing works if one or both of the slice indices is negative. If it helps, you can think of it this way: reading the list from left to right, the first element is at index 0, the second at index 1, and so on. If you use a negative index, you count from the end of the list. For example, `my_list[-1]` is the last element, `my_list[-2]` is the second-to-last element, and so on.
  - 3 Lists are zero-based, so `my_list[0:3]` returns the first three elements of the list, starting at index 0 up to but not including index 3.

### Example 3.9. Slicing Shorthand

- ```
1 # Create a list of five elements
2 my_list = ['a', 'b', 'c', 'd', 'e']
3
4 # Slice the list to get the first three elements (shorthand)
5 first_three = my_list[:3]
6
7 # Slice the list to get the last two elements (shorthand)
8 last_two = my_list[-2:]
9
10 # Print the original list and the slices
11 print(my_list)
12 print(first_three)
13 print(last_two)
```
- 1 If the left slice index is 0, you can leave it out, and 0 is implied. So `my_list[:3]` is the same as `my_list[0:3]` from Example 3.8, "Slicing a List".
 - 2 Similarly, if the right slice index is the length of the list, you can leave it out. So `my_list[-2:]` is the same as `my_list[-2:5]` because this list has five elements.
 - 3 Note the symmetry here. In this five-element list, `my_list[:3]` returns the first 3 elements, and `my_list[-2:]` returns the last two elements. In fact, `my_list[:n]` will always return the first `n` elements, and `my_list[-n:]` will always return the last `n` elements.
 - 4 If both slice indices are left out, all elements of the list are included. But this is not the same as the original list; it is a new list that happens to contain the same elements.

3.2.2. Adding Elements to Lists

Example 3.10. Adding Elements to a List

- ```
1 # Create a list of five elements
2 my_list = ['a', 'b', 'c', 'd', 'e']
3
4 # Add a single element to the end of the list
5 my_list.append('f')
6
7 # Insert a single element into a list
8 my_list.insert(2, 'g')
9
10 # Concatenate lists
11 new_list = my_list + ['h', 'i', 'j']
12
13 # Print the original list and the new lists
14 print(my_list)
15 print(new_list)
```
- 1 `append()` adds a single element to the end of the list.
  - 2 `insert()` inserts a single element into a list. The numeric argument is the index of the first element that gets bumped out of position. Note that `insert()` does not return a new list; it modifies the list in place.
  - 3 `+` concatenates lists. Note that you do not call `+` with multiple arguments; you call it with one argument, a list. In this case, that list has three elements.

### Example 3.11. The Difference between `and` and `in`

```

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

```

- 1 Lists have two methods, `append()` and `extend()` that look like they do the same thing, but are in fact completely different. `append()` takes a single argument, which is added to the end of the list.
- 2 Here you started with a list of three elements (`a`, `b`, and `c`), and you extended the list with a list of another three elements (`d`, `e`, and `f`), so you now have a list of six elements.
- 3 On the other hand, `append()` takes one argument, which can be any data type, and simply adds it to the end of the list. Here, you're calling the `append()` method on the list, passing it the list `[d, e, f]`.
- 4 Now the original list, which started as a list of three elements, contains four elements. Why four? Because the last element that you just added was the list `[d, e, f]`.

### 3.2.3. Searching Lists

#### Example 3.12. Searching a List

```

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

```

- 1 `list.index()` finds the first occurrence of a value in the list and returns the index.
- 2 `list.index()` finds the *first* occurrence of a value in the list. In this case, `0` occurs twice in the list, in `1` and `6`, but `list.index()` will return only the first index, `2`.
- 3 If the value is not found in the list, Python raises an exception. This is notably different from most languages, which will return some kind of "not found" value.
- 4 To test whether a value is in the list, use `in`, which returns `True` if the value is found or `False` if it is not.

Before version 2.2.1, Python had no separate boolean datatype. To compensate for this, Python accepted almost anything in a boolean context: `0` is false; all other numbers are true. An empty string (`''`) is false, all other strings are true.

An empty list (`[]`) is false; all other lists are true.

An empty tuple (`()`) is false; all other tuples are true.

An empty dictionary (`{}`) is false; all other dictionaries are true.

These rules still apply in Python 2.2.1 and beyond, but now you can also use an actual boolean, which has a value of `True` or `False`. Note the capitalization.

## 3.2.4. Deleting List Elements

### Example 3.13. Removing Elements from a List

```
1 # Remove elements from a list
2 my_list = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
3 # Remove the first occurrence of 'a'
4 my_list.remove('a')
5 # Remove the first occurrence of 'b'
6 my_list.remove('b')
7 # Remove the first occurrence of 'c'
8 my_list.remove('c')
9 # Remove the first occurrence of 'd'
10 my_list.remove('d')
11 # Remove the first occurrence of 'e'
12 my_list.remove('e')
13 # Remove the first occurrence of 'f'
14 my_list.remove('f')
15 # Remove the first occurrence of 'g'
16 my_list.remove('g')
17 # Remove the first occurrence of 'h'
18 my_list.remove('h')
19 # Remove the first occurrence of 'i'
20 my_list.remove('i')
21 # Remove the first occurrence of 'j'
22 my_list.remove('j')
23 # Remove the first occurrence of 'k'
24 my_list.remove('k')
25 # Remove the first occurrence of 'l'
26 my_list.remove('l')
27 # Remove the first occurrence of 'm'
28 my_list.remove('m')
29 # Remove the first occurrence of 'n'
30 my_list.remove('n')
31 # Remove the first occurrence of 'o'
32 my_list.remove('o')
33 # Remove the first occurrence of 'p'
34 my_list.remove('p')
35 # Remove the first occurrence of 'q'
36 my_list.remove('q')
37 # Remove the first occurrence of 'r'
38 my_list.remove('r')
39 # Remove the first occurrence of 's'
40 my_list.remove('s')
41 # Remove the first occurrence of 't'
42 my_list.remove('t')
43 # Remove the first occurrence of 'u'
44 my_list.remove('u')
45 # Remove the first occurrence of 'v'
46 my_list.remove('v')
47 # Remove the first occurrence of 'w'
48 my_list.remove('w')
49 # Remove the first occurrence of 'x'
50 my_list.remove('x')
51 # Remove the first occurrence of 'y'
52 my_list.remove('y')
53 # Remove the first occurrence of 'z'
54 my_list.remove('z')
55 # Print the remaining elements
56 print(my_list)
```

- 1 `remove()` removes the first occurrence of a value from a list.
- 2 `remove()` removes *only* the first occurrence of a value. In this case, `'a'` appeared twice in the list, but `remove()` removed only the first occurrence.
- 3 If the value is not found in the list, Python raises an exception. This mirrors the behavior of the `pop()` method.
- 4 `pop()` is an interesting beast. It does two things: it removes the last element of the list, and it returns the value that it removed. Note that this

## 3.2.5. Using List Operators

### Example 3.14. List Operators

```
1 # List operators
2 my_list = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
3 # Concatenate two lists
4 my_list2 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
5 # Concatenate two lists
6 my_list3 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
7 # Concatenate two lists
8 my_list4 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
9 # Concatenate two lists
10 my_list5 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
11 # Concatenate two lists
12 my_list6 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
13 # Concatenate two lists
14 my_list7 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
15 # Concatenate two lists
16 my_list8 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
17 # Concatenate two lists
18 my_list9 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
19 # Concatenate two lists
20 my_list10 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
21 # Concatenate two lists
22 my_list11 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
23 # Concatenate two lists
24 my_list12 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
25 # Concatenate two lists
26 my_list13 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
27 # Concatenate two lists
28 my_list14 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
29 # Concatenate two lists
30 my_list15 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
31 # Concatenate two lists
32 my_list16 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
33 # Concatenate two lists
34 my_list17 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
35 # Concatenate two lists
36 my_list18 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
37 # Concatenate two lists
38 my_list19 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
39 # Concatenate two lists
40 my_list20 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
41 # Concatenate two lists
42 my_list21 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
43 # Concatenate two lists
44 my_list22 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
45 # Concatenate two lists
46 my_list23 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
47 # Concatenate two lists
48 my_list24 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
49 # Concatenate two lists
50 my_list25 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
51 # Concatenate two lists
52 my_list26 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
53 # Concatenate two lists
54 my_list27 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
55 # Concatenate two lists
56 my_list28 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
57 # Concatenate two lists
58 my_list29 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
59 # Concatenate two lists
60 my_list30 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
61 # Concatenate two lists
62 my_list31 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
63 # Concatenate two lists
64 my_list32 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
65 # Concatenate two lists
66 my_list33 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
67 # Concatenate two lists
68 my_list34 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
69 # Concatenate two lists
70 my_list35 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
71 # Concatenate two lists
72 my_list36 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
73 # Concatenate two lists
74 my_list37 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
75 # Concatenate two lists
76 my_list38 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
77 # Concatenate two lists
78 my_list39 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
79 # Concatenate two lists
80 my_list40 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
81 # Concatenate two lists
82 my_list41 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
83 # Concatenate two lists
84 my_list42 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
85 # Concatenate two lists
86 my_list43 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
87 # Concatenate two lists
88 my_list44 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
89 # Concatenate two lists
90 my_list45 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
91 # Concatenate two lists
92 my_list46 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
93 # Concatenate two lists
94 my_list47 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
95 # Concatenate two lists
96 my_list48 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
97 # Concatenate two lists
98 my_list49 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
99 # Concatenate two lists
100 my_list50 = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
```

- 1 Lists can also be concatenated with the `+` operator. `my_list + my_list2` has the same result as `my_list + my_list2`. But the `+` operator returns a new (concatenated) list as a value.
- 2 Python supports the `+` operator. `my_list + my_list2` is equivalent to `my_list + my_list2`. The `+` operator works for lists, strings, and integers, and it can be overloaded to work for other types.
- 3 The `*` operator works on lists as a repeater. `my_list * 3` is equivalent to `my_list + my_list + my_list`, which concatenates the three lists into one.

### Further Reading on Lists

*How to Think Like a Computer Scientist* (<http://www.ibiblio.org/obp/thinkCSpy/>) teaches about lists and makes an important point about



*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) shows how to use lists as stacks and queues (<http://www.python.org/doc/current/faq/tuple.html>).  
Python Knowledge Base (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>) answers common questions about lists (<http://www.python.org/doc/current/faq/tuple.html>).  
*Python Library Reference* (<http://www.python.org/doc/current/lib/>) summarizes all the list methods (<http://www.python.org/doc/current/faq/tuple.html>).

## 3.3. Introducing Tuples

A tuple is an immutable list. A tuple can not be changed in any way once it is created.

### Example 3.15. Defining a tuple



- 1 A tuple is defined in the same way as a list, except that the whole set of elements is enclosed in parentheses instead of square brackets.
- 2 The elements of a tuple have a defined order, just like a list. Tuples indices are zero-based, just like a list, so the first element of a non-empty tuple is at index 0.
- 3 Negative indices count from the end of the tuple, just as with a list.
- 4 Slicing works too, just like a list. Note that when you slice a list, you get a new list; when you slice a tuple, you get a new tuple.

### Example 3.16. Tuples Have No Methods



- 1 You can't add elements to a tuple. Tuples have no `append()` method.
- 2 You can't remove elements from a tuple. Tuples have no `remove()` method.
- 3 You can't find elements in a tuple. Tuples have no `find()` method.
- 4 You can, however, use `in` to see if an element exists in the tuple.

So what are tuples good for?

Tuples are faster than lists. If you're defining a constant set of values and all you're ever going to do with it is iterate through it, use a tuple.

It makes your code safer if you "write-protect" data that does not need to be changed. Using a tuple instead of a list is like having an immutable list. Remember that I said that dictionary keys can be integers, strings, and "a few other types"? Tuples are one of those types. Tuples can be used in string formatting, as you'll see shortly.

Tuples can be converted into lists, and vice-versa. The built-in `list()` function takes a tuple and returns a list with the same elements, and the `tuple()` function takes a list and returns a tuple with the same elements.

### Further Reading on Tuples

*How to Think Like a Computer Scientist* (<http://www.ibiblio.org/obp/thinkCSpy/>) teaches about tuples and shows how to concatenate tuples. Python Knowledge Base (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>) shows how to sort a tuple (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>). *Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) shows how to define a tuple with one element (<http://www.python.org/doc/current/tut/tut.html>).

## 3.4. Declaring variables

Now that you know something about dictionaries, tuples, and lists (oh my!), let's get back to the sample program from Chapter 2, `ch2.py`.

Python has local and global variables like most other languages, but it has no explicit variable declarations. Variables spring into existence by being assigned a value.

### Example 3.17. Defining the `__main__` variable

```
if __name__ == '__main__':
```

Notice the indentation. An `if` statement is a code block and needs to be indented just like a function.

Also notice that the variable assignment is one command split over several lines, with a backslash (`"\"`) serving as a line-continuation marker.

When a command is split among several lines with the line-continuation marker (`"\"`), the continued lines can be indented in any manner; Python does not care.

Strictly speaking, expressions in parentheses, straight brackets, or curly braces (like defining a dictionary) can be split into multiple lines with the line-continuation marker.

Third, you never declared the variable `__main__` you just assigned a value to it. This is like VBScript without the `Dim` option. Luckily, unlike VBScript, Python does not require variable declarations.

### 3.4.1. Referencing Variables

#### Example 3.18. Referencing an Unbound Variable

```
x = 1
```

You will thank Python for this one day.

### 3.4.2. Assigning Multiple Values at Once

One of the cooler programming shortcuts in Python is using sequences to assign multiple values at once.

#### Example 3.19. Assigning multiple values at once

```
1 a, b, c = 'foo', 'bar', 'baz'
2
```

1 `a, b, c` is a tuple of three elements, and `a, b, c` is a tuple of three variables. Assigning one to the other assigns each of the values of `a, b, c` to each of the variables `a, b, c`.

This has all sorts of uses. I often want to assign names to a range of values. In C, you would use `enum` and manually list each constant and its associated value.

#### Example 3.20. Assigning Consecutive Values

```
1 from itertools import count
2 a, b, c = count()
3
```

- 1 The built-in `count` function returns a list of integers. In its simplest form, it takes an upper limit and returns a zero-based list counting up to that limit.
- 2 `a, b, c` and `a, b, c` are the variables you're defining. (This example came from the `calendar` module, a fun little module that prints calendars, like the one in the next example.)
- 3 Now each variable has its value: `a` is 0, `b` is 1, and so forth.

You can also use multi-variable assignment to build functions that return multiple values, simply by returning a tuple of all the values. The calendar module does this.

#### Further Reading on Variables

*Python Reference Manual* (<http://www.python.org/doc/current/ref/>) shows examples of when you can skip the line continuation character.

*How to Think Like a Computer Scientist* (<http://www.ibiblio.org/obp/thinkCSpy/>) shows how to use multi-variable assignment to swap variables.

## 3.5. Formatting Strings

Python supports formatting values into strings. Although this can include very complicated expressions, the most basic usage is to insert values into strings.

String formatting in Python uses the same syntax as the `printf` function in C.

#### Example 3.21. Introducing String Formatting

```
1
```



## 3.6. Mapping Lists

One of the most powerful features of Python is the list comprehension, which provides a compact way of mapping a list into another list by a

### Example 3.24. Introducing List Comprehensions

```
1 # Create a list of squares
2 squares = []
3 for i in range(10):
4 squares.append(i**2)
5
6 # Create a list of squares using list comprehension
7 squares = [i**2 for i in range(10)]
```

- 1 To make sense of this, look at it from right to left. `i` is the list you're mapping. Python loops through `one` element at a time, temporarily.
- 2 Note that list comprehensions do not change the original list.
- 3 It is safe to assign the result of a list comprehension to the variable that you're mapping. Python constructs the new list in memory, and

Here are the list comprehensions in the `dict` function that you declared in Chapter 2:

```
1 def dict():
2 return {k: v for k, v in data.items()}
```

First, notice that you're calling the `dict` function of the `dict` dictionary. This function returns a list of tuples of all the data in the dictionary.

### Example 3.25. The `keys` and `values` Functions

```
1 # Create a dictionary
2 data = {'a': 1, 'b': 2, 'c': 3}
3
4 # Get the keys
5 keys = data.keys()
6
7 # Get the values
8 values = data.values()
9
10 # Get the items
11 items = data.items()
```

- 1 The `keys` method of a dictionary returns a list of all the keys. The list is not in the order in which the dictionary was defined (remember that dictionaries are unordered).
- 2 The `values` method returns a list of all the values. The list is in the same order as the list returned by `keys`.
- 3 The `items` method returns a list of tuples of the form `(key, value)`. The list contains all the data in the dictionary.

Now let's see what `map` does. It takes a list, `data`, and maps it to a new list by applying string formatting to each element. The new list will have the same

### Example 3.26. List Comprehensions in `map` Step by Step

```
1 # Create a list of squares
2 squares = []
3 for i in range(10):
4 squares.append(i**2)
5
6 # Create a list of squares using list comprehension
7 squares = [i**2 for i in range(10)]
8
9 # Create a list of squares using map
10 squares = map(lambda x: x**2, range(10))
```

- 1 Note that you're using two variables to iterate through the `list`. This is another use of multi-variable assignment. The first element of `list` is `key`, and the second is `value`.
- 2 Here you're doing the same thing, but ignoring the value of `value`, so this list comprehension ends up being equivalent to `[key for key, value in list]`.
- 3 Combining the previous two examples with some simple string formatting, you get a list of strings that include both the key and value: `['key: value', 'key: value', 'key: value']`.

### Further Reading on List Comprehensions

*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) discusses another way to map lists using the built-in `map` function (<http://www.python.org/doc/current/tut/tut.html>).

*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) shows how to do nested list comprehensions (<http://www.python.org/doc/current/tut/tut.html>).

## 3.7. Joining Lists and Splitting Strings

You have a list of key-value pairs in the form `['key: value', 'key: value', 'key: value']` and you want to join them into a single string. To join any list of strings into a single string, use the `join` method.

Here is an example of joining a list from the `list` function:

```
list = ['key: value', 'key: value', 'key: value']
joined = ''.join(list)
```

One interesting note before you continue. I keep repeating that functions are objects, strings are objects... everything is an object. You might be wondering how you can call a method on a string.

The `join` method joins the elements of the list into a single string, with each element separated by a semi-colon. The delimiter doesn't need to be a string, but it must be a string-like object.

`join` works only on lists of strings; it does not do any type coercion. Joining a list that has one or more non-string elements will raise an exception.

### Example 3.27. Output of `join`

```
list = ['key: value', 'key: value', 'key: value']
joined = ''.join(list)
print(joined)
```

This string is then returned from the `join` function and printed by the calling block, which gives you the output that you marveled at when you started this chapter.

You're probably wondering if there's an analogous method to split a string into a list. And of course there is, and it's called `split`.

### Example 3.28. Splitting a String

```
joined = 'key: value key: value key: value'
list = joined.split(' ')
print(list)
```

- 1 `split` reverses `join` by splitting a string into a multi-element list. Note that the delimiter (" ") is stripped out completely; it does not appear in any of the elements of the list.
- 2 `split` takes an optional second argument, which is the number of times to split. ("Ooooooh, optional arguments..." You'll learn how to do that in a moment.)

`split` is a useful technique when you want to search a string for a substring and then work with everything before the substring (which ends up in `list[0]`).

### Further Reading on String Methods



## Chapter 4. The Power Of Introspection

This chapter covers one of Python's strengths: introspection. As you know, everything in Python is an object, and introspection is code looking at code.

### 4.1. Diving In

Here is a complete, working Python program. You should understand a good deal about it just by looking at it. The numbered lines illustrate

#### Example 4.1.

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5.0/>).

```
1 2 3
4
```

```
5
6
7
8
9
10
11
12
```

```
13 5
14
```

- 1 This module has one function, `__main__`. According to its function declaration, it takes three parameters: `obj` and `method`. The last two are actually optional.
- 2 The `__main__` function has a multi-line docstring that succinctly describes the function's purpose. Note that no return value is mentioned; this function will return `None`.
- 3 Code within the function is indented.
- 4 The `__main__` block allows this program do something useful when run by itself, without interfering with its use as a module for other programs.
- 5 Statements use `is` for comparison, and parentheses are not required.

The `__main__` function is designed to be used by you, the programmer, while working in the Python IDE. It takes any object that has functions or methods.

#### Example 4.2. Sample Usage of

```
15
16
17
18
19
20
21
22
23
24
25
26
```

By default the output is formatted to be easy to read. Multi-line strings are collapsed into a single long line, but this option can be changed by specifying `-i` on the command line.



### Example 4.3. Advanced Usage of `len`

```
def len_advanced(s):
 """Advanced len function that returns the length of a string
 and the number of vowels in the string. If the argument is not
 a string, it returns 0 and 0. If the argument is a string, it
 returns the length of the string and the number of vowels in the
 string. If the argument is a list, it returns the length of the
 list and the number of vowels in the string. If the argument is
 a tuple, it returns the length of the tuple and the number of
 vowels in the string. If the argument is a set, it returns the
 length of the set and the number of vowels in the string. If the
 argument is a dictionary, it returns the length of the dictionary
 and the number of vowels in the string. If the argument is a
 list, tuple, set, or dictionary, it returns the length of the
 object and the number of vowels in the string. If the argument
 is a string, it returns the length of the string and the
 number of vowels in the string. If the argument is not a
 string, list, tuple, set, or dictionary, it returns 0 and 0.
 """
 if not isinstance(s, (str, list, tuple, set, dict)):
 return 0, 0
 vowels = 0
 for char in s:
 if char in 'aeiou':
 vowels += 1
 return len(s), vowels
```

## 4.2. Using Optional and Named Arguments

Python allows function arguments to have default values; if the function is called without the argument, the argument gets its default value. For example, the `len` function has a default value of 0.

Here is an example of a function with two optional arguments:

```
def func(x, y=1, z=2):
 """A function with two optional arguments. x is required,
 y and z are optional, because they have default values defined.
 y is required, because it has no default value. If func is
 called with only one argument, y gets its default value of 1
 and z gets its default value of 2. If func is called with two
 arguments, y gets its default value of 1 and z gets its default
 value of 2. If func is called with three arguments, y and z
 get their specified values. In most languages, you would be out
 of luck, because you would have to specify a value for y and
 z even if you wanted to accept the default value for y and z.
 """
```

### Example 4.4. Valid Calls of `func`

- 1 With only one argument, `func(1)` gets its default value of 0 and `func(1, 1)` gets its default value of 1
- 2 With two arguments, `func(1, 2)` gets its default value of 1
- 3 Here you are naming the argument explicitly and specifying its value. `func(1, y=2)` still gets its default value of 0
- 4 Even required arguments (like `func(1, 2, 3)` which has no default value) can be named, and named arguments can appear in any order.

This looks totally whacked until you realize that arguments are simply a dictionary. The "normal" method of calling functions without arguments is just a shorthand for `func()`.

The only thing you need to do to call a function is specify a value (somehow) for each required argument; the manner and order in which you specify the arguments is up to you.

### Further Reading on Optional Arguments

*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) discusses exactly when and how default arguments are evaluated (and how to use them).

## 4.3. Using `isinstance` and Other Built-In Functions

Python has a small set of extremely useful built-in functions. All other functions are partitioned off into modules. This was actually a conscious design decision.

### 4.3.1. The `type` Function

The `type` function returns the datatype of any arbitrary object. The possible types are listed in the `sys` module. This is useful for helper functions that need to know the type of an object.

### Example 4.5. Introducing `type()`

```
1 # type() makes anything -- and I mean anything -- and returns its datatype.
2 # Integers, strings, lists, dictionaries, tuples, functions, classes, modules
3 # can take a variable and return its datatype.
4 # Also works on modules.
5 # You can use the constants in the sys module to compare types of objects.
6 # This is what the isinstance() function does, as you'll see shortly.
```

- 1 `type()` makes anything -- and I mean anything -- and returns its datatype. Integers, strings, lists, dictionaries, tuples, functions, classes, modules
- 2 `type()` can take a variable and return its datatype.
- 3 `type()` also works on modules.
- 4 You can use the constants in the `sys` module to compare types of objects. This is what the `isinstance()` function does, as you'll see shortly.

### 4.3.2. The `str()` Function

The `str()` coerces data into a string. Every datatype can be coerced into a string.

### Example 4.6. Introducing `str()`

```
1 # str() For simple datatypes like integers, you would expect str() to work,
2 # because almost every language has a function to convert an integer to a
3 # string.
4 # However, str() works on any object of any type. Here it works on a list
5 # which you've constructed in bits and pieces.
6 # Also works on modules. Note that the string representation of the module
7 # includes the pathname of the module on disk, so yours will be different.
8 # A subtle but important behavior of str() is that it works on None, the
9 # Python null value. It returns the string 'None'. You'll use this to your
10 # advantage in Chapter 5.
```

- 1 For simple datatypes like integers, you would expect `str()` to work, because almost every language has a function to convert an integer to a string.
- 2 However, `str()` works on any object of any type. Here it works on a list which you've constructed in bits and pieces.
- 3 `str()` also works on modules. Note that the string representation of the module includes the pathname of the module on disk, so yours will be different.
- 4 A subtle but important behavior of `str()` is that it works on `None`, the Python null value. It returns the string `'None'`. You'll use this to your advantage in Chapter 5.

At the heart of the `dir()` function is the powerful `__dir__` function. `__dir__` returns a list of the attributes and methods of any object: modules, functions, strings, lists, etc.

### Example 4.7. Introducing `dir()`

```
1 # dir() returns a list of the attributes and methods of any object: modules,
2 # functions, strings, lists, etc.
```



- 1 `list` is a list, so `dir(list)` returns a list of all the methods of a list. Note that the returned list contains the names of the methods as strings, not the methods themselves.
- 2 `dict` is a dictionary, so `dir(dict)` returns a list of the names of dictionary methods. At least one of these, `__str__`, should look familiar.
- 3 This is where it really gets interesting. `sys` is a module, so `dir(sys)` returns a list of all kinds of stuff defined in the module, including built-in attributes and functions.

Finally, the `callable` function takes any object and returns `True` if the object can be called, or `False` otherwise. Callable objects include functions, class methods, and some built-in functions.

### Example 4.8. Introducing `dir`



- 1 The functions in the `sys` module are deprecated (although many people still use the `sys.exit` function), but the module contains a lot of useful constants and variables.
- 2 `sys.stdout` is a function that joins a list of strings.
- 3 `sys.stdout` is not callable; it is a string. (A string does have callable methods, but the string itself is not callable.)
- 4 `sys.stdout.write` is callable; it's a function that takes two arguments.
- 5 Any callable object may have a `__call__` attribute. By using the `callable` function on each of an object's attributes, you can determine which attributes you care about.

### 4.3.3. Built-In Functions

`__builtins__` and all the rest of Python's built-in functions are grouped into a special module called `__builtins__`. (That's two underscores before and after.) If it helps, you can think of it as a module that contains all the built-in functions.

The advantage of thinking like this is that you can access all the built-in functions and attributes as a group by getting information about the `__builtins__` module.

### Example 4.9. Built-in Attributes and Functions





Python comes with excellent reference manuals, which you should peruse thoroughly to learn all the modules Python has to offer. But unlike

### Further Reading on Built-In Functions

*Python Library Reference* (<http://www.python.org/doc/current/lib/>) documents all the built-in functions (<http://www.python.org/doc/>)

## 4.4. Getting Object References With `getattr()`

You already know that Python functions are objects. What you don't know is that you can get a reference to a function without knowing its name.

### Example 4.10. Introducing `getattr()`



- 1 This gets a reference to the `pop` method of the list. Note that this is not calling the `pop` method; that would be `list.pop()`. This is the method itself.
- 2 This also returns a reference to the `pop` method, but this time, the method name is specified as a string argument to the `getattr()` function. `getattr(list, 'pop')` is an inc...
- 3 In case it hasn't sunk in just how incredibly useful this is, try this: the return value of `getattr(list, 'pop')` is the method, which you can then call just as if it were `list.pop()`.
- 4 `getattr()` also works on dictionaries.
- 5 In theory, `getattr()` would work on tuples, except that tuples have no methods, so `getattr(tuple, 'pop')` will raise an exception no matter what attribute name you give it.

### 4.4.1. `getattr()` with Modules

`getattr()` isn't just for built-in datatypes. It also works on modules.

### Example 4.11. The `getattr()` Function in `__import__`





- ❶ This returns a reference to the `function` in the `module`, which you studied in Chapter 2, *Your First Python Program*. (The hex address `0000000000000000` is just a placeholder.)
- ❷ Using `getattr` you can get the same reference to the same function. In general, `getattr(module, name)` is equivalent to `getattr(module, name)` if `module` is a module, then `name` can be anything defined in the module.
- ❸ And this is what you actually use in the `function`. `getattr` is passed into the function as an argument; `name` is a string which is the name of a method.
- ❹ In this case, `name` is the name of a function, which you can prove by getting its `__doc__`.
- ❺ Since `getattr` is a function, it is callable.

### 4.4.2. `getattr` as a Dispatcher

A common usage pattern of `getattr` is as a dispatcher. For example, if you had a program that could output data in a variety of different formats, you could use `getattr` to dispatch to the appropriate function.

For example, let's imagine a program that prints site statistics in HTML, XML, and plain text formats. The choice of output format could be based on a command-line argument or a configuration file.

#### Example 4.12. Creating a Dispatcher with `getattr`



- ❶ The `output` function takes one required argument, `format`, and one optional argument, `obj`. If `obj` is not specified, it defaults to `None` and you will end up calling `getattr(obj, format)`.
- ❷ You concatenate the `format` argument with `"output_"` to produce a function name, and then go get that function from the `module`. This allows you to dynamically get the function for the given format.
- ❸ Now you can simply call the output function in the same way as any other function. The `obj` variable is a reference to the appropriate function.

Did you see the bug in the previous example? This is a very loose coupling of strings and functions, and there is no error checking. What happens if `format` is not a valid format?

Luckily, `getattr` takes an optional third argument, a default value.

#### Example 4.13. `getattr` Default Values



❶ This function call is guaranteed to work, because you added a third argument to the call to `len`. The third argument is a default value that

As you can see, `len` is quite powerful. It is the heart of introspection, and you'll see even more powerful examples of it in later chapters.

## 4.5. Filtering Lists

As you know, Python has powerful capabilities for mapping lists into other lists, via list comprehensions (Section 3.6, “Mapping Lists”). This

Here is the list filtering syntax:

```
[expression for item in list if filter]
```

This is an extension of the list comprehensions that you know and love. The first two thirds are the same; the last part, starting with the `if` is the

### Example 4.14. Introducing List Filtering

```
1 # Create a list of numbers
2 numbers = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
3 # Filter out the even numbers
4 even_numbers = [i for i in numbers if i % 2 == 0]
5 # Print the even numbers
6 print(even_numbers)
```

❶ The mapping expression here is simple (it just returns the value of each element), so concentrate on the filter expression. As Python lo

❷ Here, you are filtering out a specific value, `2`. Note that this filters all occurrences of `2` since each time it comes up, the filter expression

❸ `count` is a list method that returns the number of times a value occurs in a list. You might think that this filter would eliminate duplicates from

Let's get back to this line from `len`:

```
len([i for i in numbers if i % 2 == 0])
```

This looks complicated, and it is complicated, but the basic structure is the same. The whole filter expression returns a list, which is assigned to

The filter expression looks scary, but it's not. You already know about `len` and `in`. As you saw in the previous section, the expression `i % 2 == 0` returns a f

So this expression takes an object (named `i`). Then it gets a list of the names of the object's attributes, methods, functions, and a few other thin

### Further Reading on Filtering Lists

*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) discusses another way to filter lists using the built-in `filter` function (<http://www.python.org/doc/current/tut/tut.html>)

## 4.6. The Peculiar Nature of `and` and `or`

In Python, `and` and `or` perform boolean logic as you would expect, but they do not return boolean values; instead, they return one of the actual values

### Example 4.15. Introducing `and`

```
1 # Create a list of numbers
2 numbers = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
3 # Filter out the even numbers
4 even_numbers = [i for i in numbers if i % 2 == 0]
5 # Print the even numbers
6 print(even_numbers)
```

c

- 1 When using `and` values are evaluated in a boolean context from left to right. `0`, `''`, `None` and `False` are false in a boolean context; everything else is true.
- 2 If any value is false in a boolean context, `and` returns the first false value. In this case, `''` is the first false value.
- 3 All values are true, so `and` returns the last value, `c`.

#### Example 4.16. Introducing `or`

1

a

2

b

3

I

4

5

6

7

a

- 1 When using `or` values are evaluated in a boolean context from left to right, just like `and`. If any value is true, `or` returns that value immediately.
- 2 `or` evaluates `''`, which is false, then `b` which is true, and returns `b`.
- 3 If all values are false, `or` returns the last value. `or` evaluates `''`, which is false, then `I` which is false, then `4` which is false, and returns `4`.
- 4 Note that `or` evaluates values only until it finds one that is true in a boolean context, and then it ignores the rest. This distinction is important.

If you're a C hacker, you are certainly familiar with the `?:` expression, which evaluates to `a` if `b` is true, and `b` otherwise. Because of the way `and` and `or`

### 4.6.1. Using the `and` Trick

#### Example 4.17. Introducing the `and` Trick

1

2

3

f

2

6

- 1 This syntax looks similar to the `?:` expression in C. The entire expression is evaluated from left to right, so the `6` is evaluated first. `6` evaluates to `True`.
- 2 `6` evaluates to `True` and then `6` evaluates to `6`.

However, since this Python expression is simply boolean logic, and not a special construct of the language, there is one extremely important

#### Example 4.18. When the `and` Trick Fails

1

2

3

6

- 1 Since `a` is an empty string, which Python considers false in a boolean context, `6` evaluates to `''`, and then `6` evaluates to `6`. Oops! That's not what we wanted.

The trick, however, will not work like the C expression `when a is false` in a boolean context.

The real trick behind the trick, then, is to make sure that the value of `a` is never false. One common way of doing this is to turn `a` into `1` and `b` into `0`.

### Example 4.19. Using the Trick Safely

```
def is_non_empty_list(x):
 return x != []
```

- 1 Since `x` is a non-empty list, it is never false. Even if `a` is `0` or some other false value, the list `x` is true because it has one element.

By now, this trick may seem like more trouble than it's worth. You could, after all, accomplish the same thing with an `if` statement, so why go through all this trouble?

### Further Reading on the Trick

Python Cookbook (<http://www.activestate.com/ASPN/Python/Cookbook/>) discusses alternatives to the trick (<http://www.activestate.com/ASPN/Python/Cookbook/recipe/10954.html>).

## 4.7. Using Lambda Functions

Python supports an interesting syntax that lets you define one-line mini-functions on the fly. Borrowed from Lisp, these so-called lambda functions can be used in a variety of ways.

### Example 4.20. Introducing Lambda Functions

```
def is_even(x):
 return x % 2 == 0

is_even(6)
True

is_even(7)
False

is_even(8)
True
```

- 1 This is a lambda function that accomplishes the same thing as the normal function above it. Note the abbreviated syntax here: there are no parentheses around the arguments, and the return value is implied by the `return` statement.
- 2 You can use a lambda function without even assigning it to a variable. This may not be the most useful thing in the world, but it just goes to show how flexible the syntax is.

To generalize, a lambda function is a function that takes any number of arguments (including optional arguments) and returns the value of a single expression.

Lambda functions are a matter of style. Using them is never required; anywhere you could use them, you could define a separate normal function and use that instead.

### 4.7.1. Real-World Lambda Functions

Here are the lambda functions in the previous example:

```
is_even = lambda x: x % 2 == 0
```

Notice that this uses the simple form of the trick, which is okay, because a lambda function is always true in a boolean context. (That doesn't mean it's always true, but it's always true in a boolean context.)



Also notice that you're using the `print` function with no arguments. You've already seen it used with one or two arguments, but without any arguments

#### Example 4.21. `print` With No Arguments

```
1 # This is a multiline string, defined by escape characters instead of triple quotes.
2 # \n is a carriage return, and \t is a tab character.
3 # print without any arguments splits on whitespace. So three spaces, a carriage return, and a tab character are all the same.
4 # You can normalize whitespace by splitting a string with .split() and then rejoining it with .join() using a single space as a delimiter. This is what the
5 # following code does.
```

- 1 This is a multiline string, defined by escape characters instead of triple quotes. `\n` is a carriage return, and `\t` is a tab character.
- 2 `print` without any arguments splits on whitespace. So three spaces, a carriage return, and a tab character are all the same.
- 3 You can normalize whitespace by splitting a string with `split()` and then rejoining it with `join()` using a single space as a delimiter. This is what the

So what is the `print` function actually doing with these `print` functions, `print`, and `print`?

`print` is now a function, but which function it is depends on the value of the `print` variable. If `print` is true, `print` will collapse whitespace; otherwise, `print` will return it

To do this in a less robust language, like Visual Basic, you would probably create a function that took a string and a `print` argument and used an `if` statement

#### Further Reading on `print` Functions

Python Knowledge Base (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>) discusses using `print` to call functions indirectly (`print func()`).  
*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) shows how to access outside variables from inside a `print` function (`print func()`).  
*The Whole Python FAQ* (<http://www.python.org/doc/FAQ.html>) has examples of obfuscated one-liners using `print` (<http://www.python.org/doc/faq/obfuscated.html>).

## 4.8. Putting It All Together

The last line of code, the only one you haven't deconstructed yet, is the one that does all the work. But by now the work is easy, because even

This is the meat of `print`

```
1 # This is a list comprehension. As you know, methods are a list of all the methods you care about in an object. So you're looping through that list with
2 # for loop.
3 # The following code shows that this is a list comprehension. As you know, methods are a list of all the methods you care about in an object. So you're looping through that list with
4 # for loop.
```

Note that this is one command, split over multiple lines, but it doesn't use the line continuation character (`\`). Remember when I said that some

Now, let's take it from the end and work backwards. The

```
1 # This is a list comprehension. As you know, methods are a list of all the methods you care about in an object. So you're looping through that list with
2 # for loop.
```

shows that this is a list comprehension. As you know, `methods` is a list of all the methods you care about in `obj`. So you're looping through that list with

#### Example 4.22. Getting a `print` Dynamically

```
1 # This is a list comprehension. As you know, methods are a list of all the methods you care about in an object. So you're looping through that list with
2 # for loop.
```



- 1 In the `__init__` function, `self` is the object you're getting help on, passed in as an argument.
- 2 As you're looping through `__dict__`, `key` is the name of the current method.
- 3 Using the `getattr` function, you're getting a reference to the `__dict__` function in the `self` module.
- 4 Now, printing the actual `__dict__` of the method is easy.

The next piece of the puzzle is the use of `str()` around the `value`. As you may recall, `str` is a built-in function that coerces data into a string. But a `self` is always

### Example 4.23. Why Use `str()` on a `self`



- 1 You can easily define a function that has no `self` so its `__dict__` attribute is `None`. Confusingly, if you evaluate the `__dict__` attribute directly, the Python IDE prints `None`.
- 2 You can verify that the value of the `__dict__` attribute is actually `None` by comparing it directly.
- 3 The `str` function takes the null value and returns a string representation of it, `None`.

In SQL, you must use `IS NULL` instead of `= NULL` to compare a null value. In Python, you can use either `is` or `==` but `is` is faster.

Now that you are guaranteed to have a string, you can pass the string to `str.ljust` which you have already defined as a function that either does or does not

Stepping back even further, you see that you're using string formatting again to concatenate the return value of `str.ljust` with the return value of `self.__dict__`.

### Example 4.24. Introducing `ljust`



- 1 `ljust` pads the string with spaces to the given length. This is what the `__init__` function uses to make two columns of output and line up all the `__dict__` in the `__dict__`.
- 2 If the given length is smaller than the length of the string, `ljust` will simply return the string unchanged. It never truncates the string.

You're almost finished. Given the padded method name from the `__dict__` method and the (possibly collapsed) `__dict__` from the call to `str`, you concatenate the

### Example 4.25. Printing a List

```
1 a = ['a', 'b', 'c']
2 print(a)
```

- ❶ This is also a useful debugging trick when you're working with lists. And in Python, you're always working with lists.

That's the last piece of the puzzle. You should now understand this code.

```
1 a = ['a', 'b', 'c']
2 print(a)
```

## 4.9. Summary

The program and its output should now make perfect sense.

```
1 a = ['a', 'b', 'c']
2 print(a)
3
4 a = ['a', 'b', 'c']
5 print(a)
6
7 a = ['a', 'b', 'c']
8 print(a)
9
10 a = ['a', 'b', 'c']
11 print(a)
```

Here is the output of

```
1 a = ['a', 'b', 'c']
2 print(a)
3
4 a = ['a', 'b', 'c']
5 print(a)
6
7 a = ['a', 'b', 'c']
8 print(a)
9
10 a = ['a', 'b', 'c']
11 print(a)
```

Before diving into the next chapter, make sure you're comfortable doing all of these things:

Defining and calling functions with optional and named arguments.

Using `str()` to coerce any arbitrary value into a string representation.

Using `getattr()` to get references to functions and other attributes dynamically.

Extending the list comprehension syntax to do list filtering.

Recognizing the `getattr()` trick and using it safely.

Defining `__call__` functions.

Assigning functions to variables and calling the function by referencing the variable. I can't emphasize this enough, because this mo

## Chapter 5. Objects and Object-Orientation

This chapter, and pretty much every chapter after this, deals with object-oriented Python programming.

### 5.1. Diving In

Here is a complete, working Python program. Read the `__doc__` of the module, the classes, and the functions to get an overview of what this program

#### Example 5.1. `__doc__`

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5.1>).

```
__doc__ = """
This module contains a complete, working Python program.
Read the __doc__ of the module, the classes, and the functions
to get an overview of what this program does.
"""

__version__ = "1.0"

__author__ = "John D. Miller"

__copyright__ = "Copyright (c) 2004 John D. Miller"

__license__ = "Python License"

__url__ = "http://diveintopython.org"

__all__ = ["__doc__", "__version__", "__author__",
 "__copyright__", "__license__", "__url__"]
```



❶ This program's output depends on the files on your hard drive. To get meaningful output, you'll need to change the directory path to p

This is the output I got on my machine. Your output will be different, unless, by some startling coincidence, you share my exact taste in musi





## 5.2. Importing Modules Using `from`

Python has two ways of importing modules. Both are useful, and you should know when to use each. One way, `import`, you've already seen in Section 5.1.

Here is the basic `from` syntax:

```
from module import *
```

This is similar to the `import` syntax that you know and love, but with an important difference: the attributes and methods of the imported module are placed directly into the current namespace.

In Python is like in Perl; in Python is like in Perl.

In Python is like In Java; In Python is like In Java.

### Example 5.2. `from sys import *`



- 1 The `sys` module contains no methods; it just has attributes for each Python object type. Note that the attribute, `sys.stdout`, must be qualified by the module name.
- 2 `sys.stdout` by itself has not been defined in this namespace; it exists only in the context of `sys`.
- 3 This syntax imports the attribute `sys.stdout` from the `sys` module directly into the local namespace.
- 4 Now `stdout` can be accessed directly, without reference to `sys`.

When should you use `from sys import *`?

If you will be accessing attributes and methods often and don't want to type the module name over and over, use `from sys import *`.

If you want to selectively import some attributes and methods but not others, use `from sys import stdout`.

If the module contains attributes or functions with the same name as ones in your module, you must use `sys.stdout` to avoid name conflicts.

Other than that, it's just a matter of style, and you will see Python code written both ways.

Use `from sys import *` sparingly, because it makes it difficult to determine where a particular function or attribute came from, and that makes debugging and re-

### Further Reading on Module Importing Techniques

eff-bot (<http://www.effbot.org/guides/>) has more to say on `from sys import *` (<http://www.effbot.org/guides/import-confusion.htm>).

*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) discusses advanced import techniques, including `from sys import *` (<http://www.pyth>).

## 5.3. Defining Classes

Python is fully object-oriented: you can define your own classes, inherit from your own or built-in classes, and instantiate the classes you've

Defined a class in Python is simple. As with functions, there is no separate interface definition. Just define the class and start coding. A Python

### Example 5.3. The Simplest Python Class



- 1 The name of this class is `Example` and it doesn't inherit from any other class. Class names are usually capitalized, but this is only a convention.
- 2 This class doesn't define any methods or attributes, but syntactically, there needs to be something in the definition, so you use `pass`. This is
- 3 You probably guessed this, but everything in a class is indented, just like the code within a function, `if` statement, `for` loop, and so forth. The



The statement in Python is like an empty set of braces ({} in Java or C).

Of course, realistically, most classes will be inherited from other classes, and they will define their own class methods and attributes. But as you

### Example 5.4. Defining the Class

```
class Animal:
 def __init__(self):
 self.name = 'Animal'
```

1 In Python, the ancestor of a class is simply listed in parentheses immediately after the class name. So the class is inherited from the class.

In Python, the ancestor of a class is simply listed in parentheses immediately after the class name. There is no special keyword like in Java.

Python supports multiple inheritance. In the parentheses following the class name, you can list as many ancestor classes as you like, separated by commas.

### 5.3.1. Initializing and Coding Classes

This example shows the initialization of the class using the method.

### Example 5.5. Initializing the Class

```
class Animal:
 def __init__(self, name):
 self.name = name
```

- 1 Classes can (and should) have too, just like modules and functions.
- 2 is called immediately after an instance of the class is created. It would be tempting but incorrect to call this the constructor of the class.
- 3 The first argument of every class method, including this always a reference to the current instance of the class. By convention, this argument is called self.
- 4 methods can take any number of arguments, and just like functions, the arguments can be defined with default values, making them optional.

By convention, the first argument of any Python class method (the reference to the current instance) is called self. This argument fills the role of this.

### Example 5.6. Coding the Class

```
class Animal:
 def __init__(self, name):
 self.name = name
 def speak(self):
 print(self.name)
```

- 1 Some pseudo-object-oriented languages like Powerbuilder have a concept of "extending" constructors and other events, where the ancestor class's constructor is called first.
- 2 I told you that this class acts like a dictionary, and here is the first sign of it. You're assigning the argument self the value of this object.
- 3 Note that the method never returns a value.

## 5.3.2. Knowing When to Use `__init__`

When defining your class methods, you *must* explicitly list `self` as the first argument for each method, including `__init__`. When you call a method of an object, Python automatically passes the object as the first argument. Whew. I realize that's a lot to absorb, but you'll get the hang of it. All Python classes work the same way, so once you learn one, you've learned them all.

Methods are optional, but when you define one, you must remember to explicitly call the ancestor's `__init__` method (if it defines one). This is more of a convention than a requirement.

### Further Reading on Python Classes

*Learning to Program* (<http://www.freenetpages.co.uk/hp/alan.gauld/>) has a gentler introduction to classes (<http://www.freenetpages.co.uk/hp/alan.gauld/python/learning-to-program/classes.html>).

*How to Think Like a Computer Scientist* (<http://www.ibiblio.org/obp/thinkCSPy/>) shows how to use classes to model compound data structures.

*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) has an in-depth look at classes, namespaces, and inheritance (<http://www.python.org/doc/current/tut/tut.html#inheritance>).

Python Knowledge Base (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>) answers common questions about classes (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>).

## 5.4. Instantiating Classes

Instantiating classes in Python is straightforward. To instantiate a class, simply call the class as if it were a function, passing the arguments that you want to use to initialize the object.

### Example 5.7. Creating a Class Instance

```
1 class Dog:
2 def __init__(self, name):
3 self.name = name
4 def bark(self):
5 print("Woof!")
```

1 You are creating an instance of the `Dog` class (defined in the `__init__.py` module) and assigning the newly created instance to the variable `fido`. You are passing the argument `"Fido"` to the `__init__` method.

2 Every class instance has a built-in attribute, `__class__`, which is the object's class. (Note that the representation of this includes the physical address of the object.)

3 You can access the instance's `__dict__` just as with a function or a module. All instances of a class share the same `__dict__`.

4 Remember when the `__init__` method assigned its `self` argument to `self`. Well, here's the result. The arguments you pass when you create the class instance are stored in the `__dict__`.

In Python, simply call a class as if it were a function to create a new instance of the class. There is no explicit `new` operator like C++ or Java.

### 5.4.1. Garbage Collection

If creating new instances is easy, destroying them is even easier. In general, there is no need to explicitly free instances, because they are freed automatically when they are no longer referenced.

### Example 5.8. Trying to Implement a Memory Leak

```
1 def leak():
2 x = []
3 while True:
4 x.append(1)
```

1 Every time the `leak` function is called, you are creating an instance of `list` and assigning it to the variable `x`, which is a local variable within the function. Since `x` is a local variable, it should be freed when the function returns. However, because `x` is a list, it is not freed, and the memory is leaked.



❶

- ❶ `__getitem__` is a normal class method; it is publicly available to be called by anyone at any time. Notice that `__getitem__` like all class methods, has `self` as its first argument.
- ❷ The `__getitem__` method of a real dictionary returns a new dictionary that is an exact duplicate of the original (all the same key-value pairs). But `MyDict` doesn't.
- ❸ You use the `__contains__` attribute to see if `MyDict` is a `dict`. If so, you're golden, because you know how to copy a `dict`: just create a new `dict` and give it the real dictionary as an argument.
- ❹ If `MyDict` is not a `dict`, then `MyDict` must be some subclass of `dict` like maybe `MyDict` in which case life gets trickier. `MyDict` doesn't know how to make an exact copy of `self`.
- ❺ The rest of the methods are straightforward, redirecting the calls to the built-in methods on `dict`.

In versions of Python prior to 2.2, you could not directly subclass built-in datatypes like strings, lists, and dictionaries. To compensate for this, Python 2.2 introduced the `__builtin__` module.

In Python, you can inherit directly from the built-in datatype, as shown in this example. There are three differences here compared to the `MyDict` example.

### Example 5.11. Inheriting Directly from Built-In Datatype

```
❶
❷
❸
❹
```

- ❶ The first difference is that you don't need to import the `__builtin__` module, since `dict` is a built-in datatype and is always available. The second is that you don't need to define `__init__`.
- ❷ The third difference is subtle but important. Because of the way `dict` works internally, it requires you to manually call its `__dict__` method to properly copy the dictionary.

#### Further Reading on `dict`

*Python Library Reference* (<http://www.python.org/doc/current/lib/>) documents the `dict` module (<http://www.python.org/doc/current/lib/>).

## 5.6. Special Class Methods

In addition to normal class methods, there are a number of special methods that Python classes can define. Instead of being called directly by the user, these methods are called by the Python interpreter.

As you saw in the previous section, normal methods go a long way towards wrapping a dictionary in a class. But normal methods alone are not enough.

### 5.6.1. Getting and Setting Items

#### Example 5.12. The `__getitem__` Special Method

```
❶
❷
❸
❹
❺
❻
❼
❽
```

- ❶ The `__getitem__` special method looks simple enough. Like the normal methods `__len__` and `__str__`, it just redirects to the dictionary to return its value. But here's the catch: `__getitem__` is called with `self` and `key`.
- ❷ This looks just like the syntax you would use to get a dictionary value, and in fact it returns the value you would expect. But here's the catch: `__getitem__` is called with `self` and `key`.

Of course, Python has a `__setitem__` special method to go along with `__getitem__` as shown in the next example.

#### Example 5.13. The `__setitem__` Special Method



- 1 Like the `__getitem__` method, `__getattr__` simply redirects to the real dictionary `__dict__` to do its work. And like `__getitem__`, you wouldn't ordinarily call it directly like this; Python does it for you.
- 2 This looks like regular dictionary syntax, except of course that `__getattr__` is really a class that's trying very hard to masquerade as a dictionary, and `__getattr__` is a special class method because it gets called for you, but it's still a class method. Just as easily as the `__getitem__` method was defined in `MP3File`, you can redefine `__getattr__` in `MP3File`.

This concept is the basis of the entire framework you're studying in this chapter. Each file type can have a handler class that knows how to get the data from the file.

For example, `MP3File` is a descendant of `File`. When an `MP3File` is set, it doesn't just set the `__dict__` (like the ancestor `File` does); it also looks in the file itself for MP3 data.

### Example 5.14. Overriding `__getattr__` in `MP3File`



- 1 Notice that this `__getattr__` method is defined exactly the same way as the ancestor method. This is important, since Python will be calling the method in `MP3File` instead of `File`.
- 2 Here's the crux of the entire `MP3File` class: if you're assigning a value to the `__dict__`, you want to do something extra.
- 3 The extra processing you do for `__dict__` is encapsulated in the `__setitem__` method. This is another class method defined in `MP3File` and when you call it, you qualify the name with `MP3File.__setitem__(self, key, value)`.
- 4 After doing this extra processing, you want to call the ancestor method. Remember that this is never done for you in Python; you must do it yourself.

When accessing data attributes within a class, you need to qualify the attribute name: `MP3File.__dict__`. When calling other methods within a class, you need to qualify the method name: `MP3File.__setitem__(self, key, value)`.

### Example 5.15. Setting an `MP3File` instance



- 1 First, you create an instance of `MP3File` without passing it a filename. (You can get away with this because the `filename` argument of the `__init__` method is optional.)

- ② Now the real fun begins. Setting the `key` of `triggers` the `__method__` on `__not__` which notices that you're setting the `key` with a real value.
- ③ Modifying the `key` will go through the same process again: Python calls `__` which calls `__` which sets all the other keys.

## 5.7. Advanced Special Class Methods

Python has more special methods than just `__` and `__`. Some of them let you emulate functionality that you may not even know about.

This example shows some of the other special methods in `__`.

### Example 5.16. More Special Methods in `__`

```

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

```

- ① `__` is a special method that is called when you call `__`. The `__` function is a built-in function that returns a string representation of an object. It
- ② `__` is called when you compare class instances. In general, you can compare any two Python objects, not just class instances, by using `==`.
- ③ `__` is called when you call `__`. The `__` function is a built-in function that returns the length of an object. It works on any object that could reasonably
- ④ `__` is called when you call `__` which you may remember as the way to delete individual items from a dictionary. When you use `del` on a class

In Java, you determine whether two string variables reference the same physical memory location by using `==`. This is called *object identity*, and

At this point, you may be thinking, "All this work just to do something in a class that I can do with a built-in datatype." And it's true that life

Special methods mean that *any class* can store key/value pairs like a dictionary, just by defining the `__` method. *Any class* can act like a sequence

While other object-oriented languages only let you define the physical model of an object ("this object has a `__` method"), Python's special class

Python has a lot of other special methods. There's a whole set of them that let classes act like numbers, allowing you to add, subtract, and do

### Further Reading on Special Class Methods

*Python Reference Manual* (<http://www.python.org/doc/current/ref/>) documents all the special class methods (<http://www.python.org>)

## 5.8. Introducing Class Attributes

You already know about data attributes, which are variables owned by a specific instance of a class. Python also supports class attributes, which

### Example 5.17. Introducing Class Attributes

```

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

```



- 1 is the class itself, not any particular instance of the class.
- 2 is a class attribute: literally, an attribute of the class. It is available before creating any instances of the class.
- 3 Class attributes are available both through direct reference to the class and through any instance of the class.

In Java, both static variables (called class attributes in Python) and instance variables (called data attributes in Python) are defined immediately.

Class attributes can be used as class-level constants (which is how you use them in Java) but they are not really constants. You can also change them.

There are no constants in Python. Everything can be changed if you try hard enough. This fits with one of the core principles of Python: bad

### Example 5.18. Modifying Class Attributes



1  
2  
2  
2

- 1 is a class attribute of the class.
- 2 is a built-in attribute of every class instance (of every class). It is a reference to the class that is an instance of (in this case, the class).
- 3 Because is a class attribute, it is available through direct reference to the class, before you have created any instances of the class.
- 4 Creating an instance of the class calls the method, which increments the class attribute by 1. This affects the class itself, not just the new instance.
- 5 Creating a second instance will increment the class attribute again. Notice how the class attribute is shared by the class and all instances.

## 5.9. Private Functions

Like most languages, Python has the concept of private elements:

- Private functions, which can't be called from outside their module.
- Private class methods, which can't be called from outside their class.
- Private attributes, which can't be accessed from outside their class.

Unlike in most languages, whether a Python function, method, or attribute is private or public is determined entirely by its name.

If the name of a Python function, class method, or attribute starts with (but doesn't end with) two underscores, it's private; everything else is public.

In this chapter, there are two methods: `__init__` and `__del__`. As you have already discussed, `__init__` is a special method; normally, you would call it indirectly by using the dictionary `__dict__`.

In Python, all special methods (like `__init__`) and built-in attributes (like `__dict__`) follow a standard naming convention: they both start with and end with two underscores.

### Example 5.19. Trying to Call a Private Method

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99  
100

- 1 If you try to call a private method, Python will raise a slightly misleading exception, saying that the method does not exist. Of course it does exist, but it's private.

### Further Reading on Private Functions

*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) discusses the inner workings of private variables (<http://www.pyth>

## 5.10. Summary

That's it for the hard-core object trickery. You'll see a real-world application of special class methods in Chapter 12, which uses `__new__` to create a new class.

The next chapter will continue using this code sample to explore other Python concepts, such as exceptions, file objects, and `__del__`.

Before diving into the next chapter, make sure you're comfortable doing all of these things:



Importing modules using either `import` or `from`

Defining and instantiating classes

Defining `__init__` methods and other special class methods, and understanding when they are called

Subclassing `dict` to define classes that act like dictionaries

Defining data attributes and class attributes, and understanding the differences between them

Defining private attributes and methods

# Chapter 6. Exceptions and File Handling

In this chapter, you will dive into exceptions, file objects, loops, and the `os` and `sys` modules. If you’ve used exceptions in another programming language, you’ll find them familiar.

## 6.1. Handling Exceptions

Like many other programming languages, Python has exception handling via `try` blocks.

Python uses `try` to handle exceptions and `raise` to generate them. Java and C++ use `try` to handle exceptions, and `throw` to generate them.

Exceptions are everywhere in Python. Virtually every module in the standard Python library uses them, and Python itself will raise them in a number of cases.

- Accessing a non-existent dictionary key will raise a `KeyError`.
- Searching a list for a non-existent value will raise a `ValueError`.
- Calling a non-existent method will raise an `AttributeError`.
- Referencing a non-existent variable will raise a `NameError`.
- Mixing datatypes without coercion will raise a `TypeError`.

In each of these cases, you were simply playing around in the Python IDE: an error occurred, the exception was printed (depending on your IDE, the error message might be slightly different).

An exception doesn’t need result in a complete program crash, though. Exceptions, when raised, can be *handled*. Sometimes an exception is raised, but the program continues to run.

### Example 6.1. Opening a Non-Existent File

```
1 # Attempt to open a non-existent file for reading.
2 f = open('nonexistent.txt', 'r')
3 # This will raise an exception.
4 # The following line will not be executed.
5 # f.read()
6 # The following line will not be executed.
7 # print(f.read())
8 # The following line will not be executed.
9 # f.close()
```

- Using the built-in `open` function, you can try to open a file for reading (more on `open` in the next section). But the file doesn’t exist, so this raises an exception.
- You’re trying to open the same non-existent file, but this time you’re doing it within a `try` block.
- When the `open` method raises an exception, you’re ready for it. The `except` line catches the exception and executes your own block of code, which prints the error message.
- Once an exception has been handled, processing continues normally on the first line after the `try` block. Note that this line will always print the error message.

Exceptions may seem unfriendly (after all, if you don’t catch the exception, your entire program will crash), but consider the alternative. Without exceptions, you’d have to check for errors at every step of the way.

#### 6.1.1. Using Exceptions For Other Purposes

There are a lot of other uses for exceptions besides handling actual error conditions. A common use in the standard Python library is to try to do something and catch the exception if it fails.

You can also define your own exceptions by creating a class that inherits from the built-in `Exception` class, and then raise your exceptions with the `raise` statement.

The next example demonstrates how to use an exception to support platform-specific functionality. This code comes from the `os` module, a wrapper for the operating system.

男  
女  
男  
女  
男  
女  
男  
女  
男  
女  
男  
女  
男  
女  
男  
女  
男  
女  
男  
女

- ## Further Reading on Exception Handling

*Python Library Reference* (<http://www.python.org/doc/current/lib/>) summarizes all the built-in exceptions (<http://www.python.org/d>

*Python Library Reference* (<http://www.python.org/doc/current/lib/>) documents the `urllib` module (<http://www.python.org/doc/current/lib/>)

Python has a built-in function, `open`, for opening a file on disk. `open` returns a file object, which has methods and attributes for getting information about the file.



- # Dive Into Python

### 6.2.1. Reading Files

After you open a file, the first thing you'll want to do is read from it, as shown in the next example.

### Example 6.4. Reading a File























- 1 A file object maintains state about the file it has open. The `tell()` method of a file object tells you your current position in the open file. Since the file is opened in read mode, the current position is at the beginning of the file.
- 2 The `seek()` method of a file object moves to another position in the open file. The second parameter specifies what the first one means; `0` means the beginning of the file, `1` means the current position, and `2` means the end of the file.
- 3 The `seek()` method confirms that the current file position has moved.
- 4 The `read()` method reads a specified number of bytes from the open file and returns a string with the data that was read. The optional parameter `size` specifies the number of bytes to read. If `size` is not specified, the default is to read until the end of the file.
- 5 The `seek()` method confirms that the current position has moved. If you do the math, you'll see that after reading 128 bytes, the position has moved to the end of the file.

### 6.2.2. Closing Files

Open files consume system resources, and depending on the file mode, other programs may not be able to access them. It's important to close

### Example 6.5. Closing a File

[illegible]



- 1 The `closed` attribute of a file object indicates whether the object has a file open or not. In this case, the file is still open (it's `False`).
- 2 To close a file, call the `close` method of the file object. This frees the lock (if any) that you were holding on the file, flushes buffered writes, and closes the underlying file descriptor.
- 3 The `closed` attribute confirms that the file is closed.
- 4 Just because a file is closed doesn't mean that the file object ceases to exist. The variable `f` will continue to exist until it goes out of scope.
- 5 Calling `close` on a file object whose file is already closed does *not* raise an exception; it fails silently.

### 6.2.3. Handling I/O Errors

Now you've seen enough to understand the file handling code in the sample code from the previous chapter. This example shows how to safely open and read a file.

#### Example 6.6. File Objects in a Block



- 1 Because opening and reading files is risky and may raise an exception, all of this code is wrapped in a `try` block. (Hey, isn't that standardized?)
- 2 The `open` function may raise an `OSError` (Maybe the file doesn't exist.)
- 3 The `read` method may raise an `OSError` (Maybe the file is smaller than 128 bytes.)
- 4 The `close` method may raise an `OSError` (Maybe the disk has a bad sector, or it's on a network drive and the network just went down.)
- 5 This is new: a `finally` block. Once the file has been opened successfully by the `open` function, you want to make absolutely sure that you close it, so you put the `close` call in a `finally` block.
- 6 At last, you handle your `Exception`. This could be the `OSError` raised by the call to `open` or `read`. Here, you really don't care, because all you want to do is print the error message.

### 6.2.4. Writing to Files

As you would expect, you can also write to files in much the same way that you read from them. There are two basic file modes:

```
"Append" mode will add data to the end of the file.·
"write" mode will overwrite the file.·
```

Either mode will create the file automatically if it doesn't already exist, so there's never a need for any sort of fiddly "if the log file doesn't exist, create it" code.

#### Example 6.7. Writing to Files



```

1
2
3
4
5
6
7
8

```

- 1 You start boldly by creating either the new file `or` overwrites the existing file, and opening the file for writing. (The second parameter
- 2 You can add data to the newly opened file with the `w` method of the file object returned by `p`
- 3 `is` a synonym for `p` This one-liner opens the file, reads its contents, and prints them.
- 4 You happen to know that `is` exists (since you just finished writing to it), so you can open it and append to it. (The `a` parameter means open
- 5 As you can see, both the original line you wrote and the second line you appended are now in `is` Also note that carriage returns are not i

### Further Reading on File Handling

*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) discusses reading and writing files, including how to read a file on  
 eff-bot (<http://www.effbot.org/guides/>) discusses efficiency and performance of various ways of reading a file (<http://www.effbot.org>)  
 Python Knowledge Base (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>) answers common questions about files (<http://www.pytho>)  
*Python Library Reference* (<http://www.python.org/doc/current/lib/>) summarizes all the file object methods (<http://www.python.org/c>)

## 6.3. Iterating with Loops

Like most other languages, Python has loops. The only reason you haven't seen them until now is that Python is good at so many other things

Most other languages don't have a powerful list datatype like Python, so you end up doing a lot of manual work, specifying a start, end, and s

### Example 6.8. Introducing the Loop

```

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

```

- 1 The syntax for a loop is similar to list comprehensions. `is` is a list, and `will` take the value of each element in turn, starting from the first
- 2 Like an `statement` or any other indented block, a loop can have any number of lines of code in it.
- 3 This is the reason you haven't seen the loop yet: you haven't needed it yet. It's amazing how often you use loops in other languages

Doing a "normal" (by Visual Basic standards) counter loop is also simple.

### Example 6.9. Simple Counters

```

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

```

```
0
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
```

- ❶ As you saw in Example 3.20, “Assigning Consecutive Values”, `range()` produces a list of integers, which you then loop through. I know it looks like you’re looping through a range, but you’re actually looping through a list of integers.
  - ❷ Don’t ever do this. This is Visual Basic-style thinking. Break out of it. Just iterate through the list, as shown in the previous example.
- Loops are not just for simple counters. They can iterate through all kinds of things. Here is an example of using a `for` loop to iterate through a dictionary.

### Example 6.10. Iterating Through a Dictionary

```
1 # Loop through the environment variables dictionary
2
3 # Print the keys and values
4
5 # Loop through the dictionary items
6
7 # Print each key and value
8
9 # Loop through the dictionary items
10
11 # Print each key and value
12
13 # Loop through the dictionary items
14
15 # Print each key and value
16
17 # Loop through the dictionary items
18
19 # Print each key and value
20
21 # Loop through the dictionary items
22
23 # Print each key and value
24
25 # Loop through the dictionary items
26
27 # Print each key and value
28
29 # Loop through the dictionary items
30
31 # Print each key and value
32
33 # Loop through the dictionary items
34
35 # Print each key and value
36
37 # Loop through the dictionary items
38
39 # Print each key and value
40
41 # Loop through the dictionary items
42
43 # Print each key and value
44
45 # Loop through the dictionary items
46
47 # Print each key and value
48
49 # Loop through the dictionary items
50
51 # Print each key and value
52
53 # Loop through the dictionary items
54
55 # Print each key and value
56
57 # Loop through the dictionary items
58
59 # Print each key and value
60
61 # Loop through the dictionary items
62
63 # Print each key and value
64
65 # Loop through the dictionary items
66
67 # Print each key and value
68
69 # Loop through the dictionary items
70
71 # Print each key and value
72
73 # Loop through the dictionary items
74
75 # Print each key and value
76
77 # Loop through the dictionary items
78
79 # Print each key and value
80
81 # Loop through the dictionary items
82
83 # Print each key and value
84
85 # Loop through the dictionary items
86
87 # Print each key and value
88
89 # Loop through the dictionary items
90
91 # Print each key and value
92
93 # Loop through the dictionary items
94
95 # Print each key and value
96
97 # Loop through the dictionary items
98
99 # Print each key and value
```

- ❶ `os.environ` is a dictionary of the environment variables defined on your system. In Windows, these are your user and system variables accessible to the current process.
- ❷ `os.environ.items()` returns a list of tuples: `(key, value)`. The `for` loop iterates through this list. The first round, it assigns `key` and `value` to `k` and `v`. In the second round, it assigns `key` and `value` to `k` and `v`.
- ❸ With multi-variable assignment and list comprehensions, you can replace the entire `for` loop with a single statement. Whether you actually do this is up to you.

Now we can look at the `for` loop in `main.py` from the sample program introduced in Chapter 5.

### Example 6.11. Loop in `main.py`

```
1 # Loop through the environment variables dictionary
2
3 # Print the keys and values
4
5 # Loop through the dictionary items
6
7 # Print each key and value
8
9 # Loop through the dictionary items
10
11 # Print each key and value
12
13 # Loop through the dictionary items
14
15 # Print each key and value
16
17 # Loop through the dictionary items
18
19 # Print each key and value
20
21 # Loop through the dictionary items
22
23 # Print each key and value
24
25 # Loop through the dictionary items
26
27 # Print each key and value
28
29 # Loop through the dictionary items
30
31 # Print each key and value
32
33 # Loop through the dictionary items
34
35 # Print each key and value
36
37 # Loop through the dictionary items
38
39 # Print each key and value
40
41 # Loop through the dictionary items
42
43 # Print each key and value
44
45 # Loop through the dictionary items
46
47 # Print each key and value
48
49 # Loop through the dictionary items
50
51 # Print each key and value
52
53 # Loop through the dictionary items
54
55 # Print each key and value
56
57 # Loop through the dictionary items
58
59 # Print each key and value
60
61 # Loop through the dictionary items
62
63 # Print each key and value
64
65 # Loop through the dictionary items
66
67 # Print each key and value
68
69 # Loop through the dictionary items
70
71 # Print each key and value
72
73 # Loop through the dictionary items
74
75 # Print each key and value
76
77 # Loop through the dictionary items
78
79 # Print each key and value
80
81 # Loop through the dictionary items
82
83 # Print each key and value
84
85 # Loop through the dictionary items
86
87 # Print each key and value
88
89 # Loop through the dictionary items
90
91 # Print each key and value
92
93 # Loop through the dictionary items
94
95 # Print each key and value
96
97 # Loop through the dictionary items
98
99 # Print each key and value
```

- ❶ It's a class attribute that defines the tags you're looking for in an MP3 file. Tags are stored in fixed-length fields. Once you read the last
- ❷ This looks complicated, but it's not. The structure of the `variables` matches the structure of the elements of the list returned by `MPRemainder`.
- ❸ Now that you've extracted all the parameters for a single MP3 tag, saving the tag data is easy. You slice from `to` to `do` to get the actual data.

## 6.4. Using `File`

Modules, like everything else in Python, are objects. Once imported, you can always get a reference to a module through the global dictionary

### Example 6.12. Introducing $\mathbb{H}$

- 1 The `sys` module contains system-level information, such as the version of Python you're running (for `sys.version`) and system-level options such as the path to the standard library (for `sys.path`).
- 2 `__import__` is a dictionary containing all the modules that have ever been imported since Python was started; the key is the module name, the value is the module object.

This example demonstrates how to use 

### Example 6.13. Using $\mathbb{R}^n$

# Dive Into Python





- 1 As new modules are imported, they are added to `__dict__`. This explains why importing the same module twice is very fast: Python has already imported it.
- 2 Given the name (as a string) of any previously-imported module, you can get a reference to the module itself through the `__dict__`.

The next example shows how to use the `__class__` attribute with the `__dict__` to get a reference to the module in which a class is defined.

### Example 6.14. The `__class__` Attribute



- 1 Every Python class has a built-in class attribute `__module__` which is the name of the module in which the class is defined.
- 2 Combining this with the `__dict__`, you can get a reference to the module in which a class is defined.

Now you're ready to see how `__module__` is used in `__main__` the sample program introduced in Chapter 5. This example shows that portion of the code.

### Example 6.15. `__main__`



- 1 This is a function with two arguments; `__name__` is required, but `__module__` is optional and defaults to the module that contains the `__main__` class. This looks inefficient.
- 2 You'll plow through this line later, after you dive into the `__main__` module. For now, take it on faith that `__main__` ends up as the name of a class, like `__main__`.
- 3 You already know about `__dict__` which gets a reference to an object by name. `__module__` is a complementary function that checks whether an object has a `__module__` attribute.

### Further Reading on Modules

*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) discusses exactly when and how default arguments are evaluated (`__name__` and `__module__`).

*Python Library Reference* (<http://www.python.org/doc/current/lib/>) documents the `__dict__` (<http://www.python.org/doc/current/lib/module-dict.html>).

## 6.5. Working with Directories

The `os` module has several functions for manipulating files and directories. Here, we're looking at handling pathnames and listing the contents of a directory.

### Example 6.16. Constructing Pathnames





- 1 is a reference to a module -- which module depends on your platform. Just as encapsulates differences between platforms by setting
- 2 The function of constructs a pathname out of one or more partial pathnames. In this case, it simply concatenates strings. (Note that de
- 3 In this slightly less trivial case, will add an extra backslash to the pathname before joining it to the filename. I was overjoyed when I c
- 4 will expand a pathname that uses to represent the current user's home directory. This works on any platform where users have a hom
- 5 Combining these techniques, you can easily construct pathnames for directories and files under the user's home directory.

### Example 6.17. Splitting Pathnames



- 1 The function splits a full pathname and returns a tuple containing the path and filename. Remember when I said you could use multi-
- 2 You assign the return value of the function into a tuple of two variables. Each variable receives the value of the corresponding elemen
- 3 The first variable, receives the value of the first element of the tuple returned from the file path.
- 4 The second variable, receives the value of the second element of the tuple returned from the filename.
- 5 Also contains a function which splits a filename and returns a tuple containing the filename and the file extension. You use the same

### Example 6.18. Listing Directories





- 1 The `os.listdir()` function takes a pathname and returns a list of the contents of the directory.
- 2 It returns both files and folders, with no indication of which is which.
- 3 You can use list filtering and the `os.path.isfile()` function of the `os.path` module to separate the files from the folders. `os.path.isfile()` takes a pathname and returns 1 if the path is a file, and 0 otherwise.
- 4 `os.path.isdir()` also has a function which returns 1 if the path represents a directory, and 0 otherwise. You can use this to get a list of the subdirectories.

### Example 6.19. Listing Directories in `os`



- 1 `os.listdir()` returns a list of all the files and folders in `dir`.
- 2 Iterating through the list with `for`, you use `os.path.normcase()` to normalize the case according to operating system defaults. `os.path.normcase()` is a useful little function that converts the pathname to the case of the operating system.
- 3 Iterating through the normalized list with `for` again, you use `os.path.splitext()` to split each filename into name and extension.
- 4 For each file, you see if the extension is in the list of file extensions you care about (which was passed to the `main()` function).
- 5 For each file you care about, you use `os.path.join()` to construct the full pathname of the file, and return a list of the full pathnames.

Whenever possible, you should use the functions in `os.path` for file, directory, and path manipulations. These modules are wrappers for platform-specific functions.

There is one other way to get the contents of a directory. It's very powerful, and it uses the sort of wildcards that you may already be familiar with.

### Example 6.20. Listing Directories with `glob`



- 1 As you saw earlier, `ls` simply takes a directory path and lists all files and directories in that directory.
- 2 The `find` module, on the other hand, takes a wildcard and returns the full path of all files and directories matching the wildcard. Here the wildcard is `*.mp3`.
- 3 If you want to find all the files in a specific directory that start with "s" and end with ".mp3", you can do that too.
- 4 Now consider this scenario: you have a `dir` directory, with several subdirectories within it, with `files` within each subdirectory. You can get a list of all files in `dir` with the following command:

## Further Reading on the eModule

Python Knowledge Base (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>) answers questions about the `os` module (<http://www.python.org/doc/2.6.0/library/os/>).

*Python Library Reference* (<http://www.python.org/doc/current/lib/>) documents the [queue module](http://www.python.org/doc/current/lib/module-queue.html) ([http://www.python.org/doc/current/lib/module-](http://www.python.org/doc/current/lib/module-queue.html)

## 6.6. Putting It All Together

Once again, all the dominoes are in place. You've seen how each line of code works. Now let's step back and see how it all fits together.

### Example 6.21.

- 1 Is the main attraction of this entire module. It takes a directory (like `in` in my case) and a list of interesting file extensions (like `['.py']` and it returns a list of the full pathnames of all the files in that directory that have an interesting file extension.
- 2 As you saw in the previous section, this line of code gets a list of the full pathnames of all the files in `in` that have an interesting file extension.
- 3 Old-school Pascal programmers may be familiar with them, but most people give me a blank stare when I tell them that Python supports them.
- 4 Now that you've seen the `os.path.splitext` module, this line should make more sense. It gets the extension of the file (`f` forces it to uppercase (`f.upper()` slices off the extension).
- 5 Having constructed the name of the handler class that would handle this file, you check to see if that handler class actually exists in the `handlers` module.
- 6 For each file in the "interesting files" list (`files`) you call `handlers.get_handler` with the filename (`f`). Calling `handlers.get_handler` returns a class; you don't know exactly which class

Note that `File` is completely generic. It doesn't know ahead of time which types of files it will be getting, or which classes are defined that could

## 6.7. Summary

The program introduced in Chapter 5 should now make perfect sense.

[illegible]



Before diving into the next chapter, make sure you're comfortable doing the following things:

- Catching exceptions with `try/except`

- Protecting external resources with `with`

- Reading from files

- Assigning multiple values at once in a `for` loop

- Using the `os` module for all your cross-platform file manipulation needs

- Dynamically instantiating classes of unknown type by treating classes as objects and passing them around

## Chapter 7. Regular Expressions

Regular expressions are a powerful and standardized way of searching, replacing, and parsing text with complex patterns of characters. If you

## 7.1. Diving In

Strings have methods for searching (`find` and `findall`), replacing (`replace`) and parsing (`parse`) but they are limited to the simplest of cases. The search methods look

If what you're trying to do can be accomplished with string functions, you should use them. They're fast and simple and easy to read, and the

Although the regular expression syntax is tight and unlike normal code, the result can end up being *more* readable than a hand-rolled solution

## 7.2. Case Study: Street Addresses

This series of examples was inspired by a real-life problem I had in my day job several years ago, when I needed to scrub and standardize str

### Example 7.1. Matching at the End of a String

| Mode of Transport | Number of people |
|-------------------|------------------|
| Car               | 4                |
| Train             | 3                |
| Bus               | 2                |
| Bicycle           | 1                |
| On foot           | 1                |

- 1 My goal is to standardize a street address so that `As` always abbreviated as `A`. At first glance, I thought this was simple enough that I could write a function like this:
- 2 Life, unfortunately, is full of counterexamples, and I quickly discovered this one. The problem here is that `A` appears twice in the address `1000 Avenue of the Americas`.
- 3 To solve the problem of addresses with more than one `A` substring, you could resort to something like this: only search and replace `A` on the last occurrence of `A` in the string.
- 4 It's time to move up to regular expressions. In Python, all functionality related to regular expressions is contained in the `re` module.
- 5 Take a look at the first parameter: `A$`. This is a simple regular expression that matches `A` only when it occurs at the end of a string. The `$` character is the end-of-string anchor.
- 6 Using the `sub` function, you search the string `s` for the regular expression `A$` and replace it with `A`. This matches the `A` at the end of the string `Sub`.

Continuing with my story of scrubbing addresses, I soon discovered that the previous example, matching `^` at the end of the address, was not good.

### Example 7.2. Matching Whole Words



- 1 What I *really* wanted was to match `\b` when it was at the end of the string *and* it was its own whole word, not a part of some larger word
- 2 To work around the backslash plague, you can use what is called a raw string, by prefixing the string with the letter `r` This tells Python
- 3 *\*sigh\** Unfortunately, I soon found more cases that contradicted my logic. In this case, the street address contained the word `\b` as a whole
- 4 To solve this problem, I removed the `\` character and added another `\b` Now the regular expression reads "match `\b` when it's a whole word"

## 7.3. Case Study: Roman Numerals

You've most likely seen Roman numerals, even if you didn't recognize them. You may have seen them in copyrights of old movies and television shows.

In Roman numerals, there are seven characters that are repeated and combined in various ways to represent numbers.

I= 1  
V= 5  
X= 10  
L= 50  
C= 100  
D= 500  
M= 1000

The following are some general rules for constructing Roman numerals:

Characters are additive. `II` is 2 and `III` is 3 `VI` literally, "5 and 1", `VII` is 7 and `VIII` is 8

The tens characters (`IX` and `M`) can be repeated up to three times. At 4 you need to subtract from the next highest fives character. You

Similarly, at 9 you need to subtract from the next highest tens character: `IX` but `IX` is less than `X` not `V` since the `I` character can not be

The fives characters can not be repeated. The number `4` is always represented as `IV` never as `IIII` The number `9` is always `IX` never `IIIIIIII`

Roman numerals are always written highest to lowest, and read left to right, so the order the of characters matters very much. `DC` is 600

### 7.3.1. Checking for Thousands

What would it take to validate that an arbitrary string is a valid Roman numeral? Let's take it one digit at a time. Since Roman numerals are a

#### Example 7.3. Checking for Thousands



- 1 This pattern has three parts:  
    `^` to match what follows only at the beginning of the string. If this were not specified, the pattern would match no matter where



`M` optionally match a single `M` character. Since this is repeated three times, you're matching anywhere from zero to three `M` characters.

`$` to match what precedes only at the end of the string. When combined with the `^` character at the beginning, this means that the

- ❷ The essence of the `re` module is the `re.match()` function, that takes a regular expression (`pattern`) and a string (`string`) to try to match against the regular expression.
- ❸ `re.match()` matches because the first and second optional `M` characters match and the third `M` is ignored.
- ❹ `re.match()` matches because all three `M` characters match.
- ❺ `re.match()` does not match. All three `M` characters match, but then the regular expression insists on the string ending (because of the `$` character), and it doesn't.
- ❻ Interestingly, an empty string also matches this regular expression, since all the `M` characters are optional.

### 7.3.2. Checking for Hundreds

The hundreds place is more difficult than the thousands, because there are several mutually exclusive ways it could be expressed, depending on

```
0- C
0- C
0- C
0- D
0- D
0- D
0- D
0- D
0- M
```

So there are four possible patterns:

```
M
D
Zero to three C characters (zero if the hundreds place is 0)
D followed by zero to three C characters
```

The last two patterns can be combined:

```
an optional D followed by zero to three C characters
```

This example shows how to validate the hundreds place of a Roman numeral.

#### Example 7.4. Checking for Hundreds

```
1
2
3
4
5
6
7
8
9
10
```

- ❶ This pattern starts out the same as the previous one, checking for the beginning of the string (`^`), then the thousands place (`M`). Then it has



- ❹ This matches the start of the string, then three ~~M~~ut of a possible three, then the end of the string.
- ❺ This matches the start of the string, then three ~~M~~ut of a possible three, but then *does not match* the end of the string. The regular expres

There is no way to programmatically determine that two regular expressions are equivalent. The best you can do is write a lot of test cases to

### 7.4.1. Checking for Tens and Ones

Now let's expand the Roman numeral regular expression to cover the tens and ones place. This example shows the check for tens.

#### Example 7.7. Checking for Tens

```

❹
❺
❻
❼
❸
❷
❶
>

```

- ❶ This matches the start of the string, then the first optional ~~M~~then ~~M~~then ~~X~~then the end of the string. Remember, the ~~M~~ syntax means "match
- ❷ This matches the start of the string, then the first optional ~~M~~then ~~M~~then ~~X~~Of the ~~X~~it matches the ~~l~~and skips all three optional ~~X~~characters.
- ❸ This matches the start of the string, then the first optional ~~M~~then ~~M~~then the optional ~~l~~and the first optional ~~X~~skips the second and third op
- ❹ This matches the start of the string, then the first optional ~~M~~then ~~M~~then the optional ~~l~~and all three optional ~~X~~characters, then the end of th
- ❺ This matches the start of the string, then the first optional ~~M~~then ~~M~~then the optional ~~l~~and all three optional ~~X~~characters, then *fails to match*

The expression for the ones place follows the same pattern. I'll spare you the details and show you the end result.

```

❹
❺
❻
❼
❸
❷
❶
>

```

So what does that look like using this alternate ~~M~~ syntax? This example shows the new syntax.

#### Example 7.8. Validating Roman Numerals with ~~M~~

```

❹
❺
❻
❼
❸
❷
❶
>

```



- ❶ This matches the start of the string, then one of a possible four **M**characters, then **0**Of that, it matches the optional **D**and zero of three pos
- ❷ This matches the start of the string, then two of a possible four **M**characters, then the **0**with a **D**and one of three possible **C**characters; then
- ❸ This matches the start of the string, then four out of four **M**characters, then **0**with a **D**and three out of three **C**characters; then **0**with an **l**and
- ❹ Watch closely. (I feel like a magician. "Watch closely, kids, I'm going to pull a rabbit out of my hat.") This matches the start of the str

If you followed all that and understood it on the first try, you're doing better than I did. Now imagine trying to understand someone else's reg

In the next section you'll explore an alternate syntax that can help keep your expressions maintainable.

## 7.5. Verbose Regular Expressions

So far you've just been dealing with what I'll call "compact" regular expressions. As you've seen, they are difficult to read, and even if you f

Python allows you to do this with something called *verbose regular expressions*. A verbose regular expression is different from a compact reg

Whitespace is ignored. Spaces, tabs, and carriage returns are not matched as spaces, tabs, and carriage returns. They're not matched  
Comments are ignored. A comment in a verbose regular expression is just like a comment in Python code: it starts with a **#**character

This will be more clear with an example. Let's revisit the compact regular expression you've been working with, and make it a verbose regul

### Example 7.9. Regular Expressions with Inline Comments



- ❶ The most important thing to remember when using verbose regular expressions is that you need to pass an extra argument when worki
- ❷ This matches the start of the string, then one of a possible four **M**then **0**then **l**and three of a possible three **X**then **X**then the end of the str
- ❸ This matches the start of the string, then four of a possible four **M**then **D**and three of a possible three **C**then **l**and three of a possible three
- ❹ This does not match. Why? Because it doesn't have the **0**flag, so the **re**function is treating the pattern as a compact regular expression, w

## 7.6. Case study: Parsing Phone Numbers

So far you’ve concentrated on matching whole patterns. Either the pattern matches, or it doesn’t. But regular expressions are much more powerful.

This example came from another real-world problem I encountered, again from a previous day job. The problem: parsing an American phone number.

Here are the phone numbers I needed to be able to accept:

```
000-000-0000
000-000-0000
000-000-0000
000-000-0000
000-000-0000
000-000-0000
000-000-0000
000-000-0000
000-000-0000
```

Quite a variety! In each of these cases, I need to know that the area code was **000**, the trunk was **000**, and the rest of the phone number was **0000**. For those who are not familiar with the format, the area code is the first three digits, the trunk is the next three digits, and the rest is the rest of the phone number.

Let’s work through developing a solution for phone number parsing. This example shows the first step.

### Example 7.10. Finding Numbers

```
>>> re.compile(r'\d{3}-\d{3}-\d{4}')
>>>
```

- 1 Always read regular expressions from left to right. This one matches the beginning of the string, and then **\d{3}**. What’s **\d{3}**? Well, the **\d** means “any digit” and the **{3}** means “exactly three times”.
- 2 To get access to the groups that the regular expression parser remembered along the way, use the **groups()** method on the object that the **match()** function returns.
- 3 This regular expression is not the final answer, because it doesn’t handle a phone number with an extension on the end. For that, you’ll need to use **re.compile(r'\d{3}-\d{3}-\d{4}-\d{4}')**.

### Example 7.11. Finding the Extension

```
>>> re.compile(r'\d{3}-\d{3}-\d{4}-\d{4}')
>>>
```

- 1 This regular expression is almost identical to the previous one. Just as before, you match the beginning of the string, then a **\d{3}**, then a **-**, then a **\d{3}**, then a **-**, then a **\d{4}**, then a **-**, then a **\d{4}**.
- 2 The **groups()** method now returns a tuple of four elements, since the regular expression now defines four groups to remember.
- 3 Unfortunately, this regular expression is not the final answer either, because it assumes that the different parts of the phone number are separated by dashes.
- 4 Oops! Not only does this regular expression not do everything you want, it’s actually a step backwards, because now you can’t parse phone numbers with parentheses around the area code.

The next example shows the regular expression to handle separators between the different parts of the phone number.

### Example 7.12. Handling Different Separators

```
>>> re.compile(r'\d{3}[-]\d{3}[-]\d{4}')
>>>
```

- 1 Hang on to your hat. You're matching the beginning of the string, then a group of three digits, then `-`. What the heck is that? Well, `-` means you can now match phone numbers where the parts are separated by spaces instead of hyphens.
- 2 Using `-` instead of `-` means you can now match phone numbers where the parts are separated by spaces instead of hyphens.
- 3 Of course, phone numbers separated by hyphens still work too.
- 4 Unfortunately, this is still not the final answer, because it assumes that there is a separator at all. What if the phone number is entered without a separator?
- 4 Oops! This still hasn't fixed the problem of requiring extensions. Now you have two problems, but you can solve both of them with the next example.

The next example shows the regular expression for handling phone numbers *without* separators.

### Example 7.13. Handling Numbers Without Separators

```
>>> re.compile(r'\d{3}\d{3}\d{4}')
>>>
```

- 1 The only change you've made since that last step is changing all the `-` to `*`. Instead of `-` between the parts of the phone number, you now match `*`.
- 2 Lo and behold, it actually works. Why? You matched the beginning of the string, then a remembered group of three digits (`\d{3}`), then zero or more digits (`\d*`).
- 3 Other variations work now too: dots instead of hyphens, and both a space and an `x` before the extension.
- 4 Finally, you've solved the other long-standing problem: extensions are optional again. If no extension is found, the `match` method still returns `None`.
- 5 I hate to be the bearer of bad news, but you're not finished yet. What's the problem here? There's an extra character before the area code.

The next example shows how to handle leading characters in phone numbers.

### Example 7.14. Handling Leading Characters

```
>>> re.compile(r'[1-9]\d{9}')
>>>
```



- 1 This is the same as in the previous example, except now you're matching zero or more non-numeric characters, before the first remainder.
- 2 You can successfully parse the phone number, even with the leading left parenthesis before the area code. (The right parenthesis after the area code is optional.)
- 3 Just a sanity check to make sure you haven't broken anything that used to work. Since the leading characters are entirely optional, this should work.
- 4 This is where regular expressions make me want to gouge my eyes out with a blunt object. Why doesn't this phone number match? Because the `+` is greedy, and it's matching the trailing `+` as well.

Let's back up for a second. So far the regular expressions have all matched from the beginning of the string. But now you see that there may

### Example 7.15. Phone Number, Wherever I May Find Ye



- 1 Note the lack of `^` in this regular expression. You are not matching the beginning of the string anymore. There's nothing that says you must match the beginning of the string.
- 2 Now you can successfully parse a phone number that includes leading characters and a leading digit, plus any number of any kind of separator.
- 3 Sanity check. this still works.
- 4 That still works too.

See how quickly a regular expression can get out of control? Take a quick glance at any of the previous iterations. Can you tell the difference

While you still understand the final answer (and it is the final answer; if you've discovered a case it doesn't handle, I don't want to know about it).

### Example 7.16. Parsing Phone Numbers (Final Version)





- ❶ Other than being spread out over multiple lines, this is exactly the same regular expression as the last step, so it's no surprise that it passes.
- ❷ Final sanity check. Yes, this still works. You're done.

### Further Reading on Regular Expressions

Regular Expression HOWTO (<http://py-howto.sourceforge.net/regex/regex.html>) teaches about regular expressions and how to use them.  
*Python Library Reference* (<http://www.python.org/doc/current/lib/>) summarizes the `re` module (<http://www.python.org/doc/current/lib/>)

## 7.7. Summary

This is just the tiniest tip of the iceberg of what regular expressions can do. In other words, even though you're completely overwhelmed by the

You should now be familiar with the following techniques:

- `^` matches the beginning of a string.
- `$` matches the end of a string.
- `\b` matches a word boundary.
- `\d` matches any numeric digit.
- `\D` matches any non-numeric character.
- `X?` matches an optional `X` character (in other words, it matches an `X` zero or one times).
- `X*` matches `X` zero or more times.
- `X+` matches `X` one or more times.
- `X{n,m}` matches an `X` character at least `n` times, but not more than `m` times.
- `X|Y` matches either `X` or `Y`.
- `(...)` in general is a *remembered group*. You can get the value of what matched by using the `group()` method of the object returned by `re.match()`.

Regular expressions are extremely powerful, but they are not the correct solution for every problem. You should learn enough about them to

Some people, when confronted with a problem, think "I know, I'll use regular expressions." Now they have two problems.  
--Jamie Zawinski, in c



## Chapter 8. HTML Processing

## 8.1. Diving in

I often see questions on comp.lang.python (<http://groups.google.com/groups?group=comp.lang.python>) like "How can I list all the [headers|

Here is a complete, working Python program in two parts. The first part, `html.py`, is a generic tool to help you process HTML files by walking through

### Example 8.1.

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5>).

1

2  
3  
4  
5  
6

7  
8  
9  
10  
11  
12  
13

14  
15  
16  
17

18  
19  
20  
21  
22  
23  
24

25  
26  
27

**Example 8.2.**

28  
29

30  
31

32  
33  
34  
35  
36

37  
38  
39  
40  
41

42  
43

```
def __init__(self, data):
 self.data = data

def __str__(self):
 return str(self.data)

def __repr__(self):
 return repr(self.data)

def __len__(self):
 return len(self.data)

def __getitem__(self, index):
 return self.data[index]

def __setitem__(self, index, value):
 self.data[index] = value

def __delitem__(self, index):
 del self.data[index]

def __iter__(self):
 return iter(self.data)

def __contains__(self, item):
 return item in self.data

def __add__(self, other):
 return self.data + other

def __sub__(self, other):
 return self.data - other

def __mul__(self, other):
 return self.data * other

def __div__(self, other):
 return self.data / other

def __mod__(self, other):
 return self.data % other

def __pow__(self, other):
 return self.data ** other

def __lt__(self, other):
 return self.data < other

def __le__(self, other):
 return self.data <= other

def __gt__(self, other):
 return self.data > other

def __ge__(self, other):
 return self.data >= other

def __eq__(self, other):
 return self.data == other

def __ne__(self, other):
 return self.data != other
```





|                        |                                                                                                                                                                                                                       |
|------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Entity reference       | An escaped character referenced by its decimal or hexadecimal equivalent, like <code>&amp;#x27;</code> . When found, <code>html.parser.HTMLParser</code> calls <code>unescape_entity()</code> with the text of the de |
| Comment                | An HTML entity, like <code>&amp;lt;!--</code> . When found, <code>html.parser.HTMLParser</code> calls <code>unescape_entity()</code> with the name of the HTML entity.                                                |
| Processing instruction | An HTML comment, enclosed in <code>&lt;!--</code> . When found, <code>html.parser.HTMLParser</code> calls <code>unescape_entity()</code> with the body of the comment.                                                |
| Declaration            | An HTML processing instruction, enclosed in <code>&lt;?&gt;</code> . When found, <code>html.parser.HTMLParser</code> calls <code>unescape_entity()</code> with the body of the processing instruction.                |
| Text data              | An HTML declaration, such as a <code>&lt;!--</code> enclosed in <code>&lt;!--</code> . When found, <code>html.parser.HTMLParser</code> calls <code>unescape_entity()</code> with the body of the declaration.         |
|                        | A block of text. Anything that doesn't fit into the other 7 categories. When found, <code>html.parser.HTMLParser</code> calls <code>unescape_entity()</code> with the text.                                           |

Python 2.0 had a bug where `HTMLParser` would not recognize declarations at all (`unescape_entity()` would never be called), which meant that `HTMLParser` were silently ignored. This comes with a test suite to illustrate this. You can run `test_htmlparser.py` passing the name of an HTML file on the command line, and it will print out the tags a

In the ActivePython IDE on Windows, you can specify command line arguments in the "Run script" dialog. Separate multiple arguments with

### Example 8.4. Sample test of `HTMLParser`

Here is a snippet from the table of contents of the HTML version of this book. Of course your paths may vary. (If you haven't downloaded th

```
<table>
<tr>
<td>Entity reference</td>
<td>An escaped character referenced by its decimal or hexadecimal equivalent, like <code>'</code>. When found, <code>html.parser.HTMLParser</code> calls <code>unescape_entity()</code> with the text of the de</td>
</tr>
<tr>
<td>Comment</td>
<td>An HTML entity, like <code><!--</code>. When found, <code>html.parser.HTMLParser</code> calls <code>unescape_entity()</code> with the name of the HTML entity.</td>
</tr>
<tr>
<td>Processing instruction</td>
<td>An HTML comment, enclosed in <code><!--</code>. When found, <code>html.parser.HTMLParser</code> calls <code>unescape_entity()</code> with the body of the comment.</td>
</tr>
<tr>
<td>Declaration</td>
<td>An HTML processing instruction, enclosed in <code><?></code>. When found, <code>html.parser.HTMLParser</code> calls <code>unescape_entity()</code> with the body of the processing instruction.</td>
</tr>
<tr>
<td>Text data</td>
<td>An HTML declaration, such as a <code><!--</code> enclosed in <code><!--</code>. When found, <code>html.parser.HTMLParser</code> calls <code>unescape_entity()</code> with the body of the declaration.</td>
</tr>
<tr>
<td></td>
<td>A block of text. Anything that doesn't fit into the other 7 categories. When found, <code>html.parser.HTMLParser</code> calls <code>unescape_entity()</code> with the text.</td>
</tr>
</table>
```

Running this through the test suite of `HTMLParser` yields this output:

```
<table>
<tr>
<td>Entity reference</td>
<td>An escaped character referenced by its decimal or hexadecimal equivalent, like <code>'</code>. When found, <code>html.parser.HTMLParser</code> calls <code>unescape_entity()</code> with the text of the de</td>
</tr>
<tr>
<td>Comment</td>
<td>An HTML entity, like <code><!--</code>. When found, <code>html.parser.HTMLParser</code> calls <code>unescape_entity()</code> with the name of the HTML entity.</td>
</tr>
<tr>
<td>Processing instruction</td>
<td>An HTML comment, enclosed in <code><!--</code>. When found, <code>html.parser.HTMLParser</code> calls <code>unescape_entity()</code> with the body of the comment.</td>
</tr>
<tr>
<td>Declaration</td>
<td>An HTML processing instruction, enclosed in <code><?></code>. When found, <code>html.parser.HTMLParser</code> calls <code>unescape_entity()</code> with the body of the processing instruction.</td>
</tr>
<tr>
<td>Text data</td>
<td>An HTML declaration, such as a <code><!--</code> enclosed in <code><!--</code>. When found, <code>html.parser.HTMLParser</code> calls <code>unescape_entity()</code> with the body of the declaration.</td>
</tr>
<tr>
<td></td>
<td>A block of text. Anything that doesn't fit into the other 7 categories. When found, <code>html.parser.HTMLParser</code> calls <code>unescape_entity()</code> with the text.</td>
</tr>
</table>
```

Here's the roadmap for the rest of the chapter:

- Subclass `HTMLParser` to create classes that extract interesting data out of HTML documents.

- Subclass `HTMLParser` to create `HTMLReconstructor` which overrides all 8 handler methods and uses them to reconstruct the original HTML from the pieces.

- Subclass `HTMLReconstructor` to create `HTMLReconstructor` which adds some methods to process specific HTML tags specially, and overrides the `unescape` method to provide a framework for processing escaped HTML.

- Subclass `HTMLReconstructor` to create classes that define text processing rules used by `HTMLReconstructor`.

- Write a test suite that grabs a real web page from `urllib` and processes it.

Along the way, you'll also learn about `re` and dictionary-based string formatting.

## 8.3. Extracting data from HTML documents

To extract data from HTML documents, subclass the `HTMLParser` class and define methods for each tag or entity you want to capture.

The first step to extracting data from an HTML document is getting some HTML. If you have some HTML lying around on your hard drive,

### Example 8.5. Introducing `urllib`



- 1 The `urllib` module is part of the standard Python library. It contains functions for getting information about and actually retrieving data from the web.
- 2 The simplest use of `urllib` is to retrieve the entire text of a web page using the `urllib.urlopen` function. Opening a URL is similar to opening a file. The returned object is a file-like object.
- 3 The simplest thing to do with the file-like object returned by `urllib.urlopen` is to read the entire HTML of the web page into a single string. This is done by calling the `read` method.
- 4 When you're done with the object, make sure to `close` it, just like a normal file object.
- 5 You now have the complete HTML of the home page of `http://www.python.org` in a string, and you're ready to parse it.

### Example 8.6. Introducing $\mathbf{P}$

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5>).

- 1 `__call__` is called by the `Parser` method of `Parser` and it can also be called manually once an instance of the parser has been created. So if you need to do an action when a tag is found, you can override `__call__`.
- 2 `__call__` is called by `Parser` whenever it finds an `etag`. The tag may contain an `attribute`, and/or other attributes, like `name` or `id`. The `params` parameter is a list of tuples containing the attribute name and value.
- 3 You can find out whether this `etag` has an `attribute` with a simple multi-variable list comprehension.
- 4 String comparisons like `name == 'h1'` are always case-sensitive, but that's safe in this case, because `Parser` converts attribute names to lowercase while building the `tree`.

### Example 8.7. Using

1. Deletes `inserting`, which is what gets HTML into the parser.
2. Like files, you should `close` your URL objects as soon as you're done with them.
3. You should `close` your parser object, too, but for a different reason. You've read all the data and fed it to the parser, but the `close` method isn't guaranteed to be there.
4. Once the parser is `done`, the parsing is complete, and `urls` contains a list of all the linked URLs in the HTML document. (Your output may look like this.)



## 8.4. Introducing `HTMLParser`

`HTMLParser` doesn't produce anything by itself. It parses and parses and parses, and it calls a method for each interesting thing it finds, but the methods are defined by subclasses. `HTMLParser` and provides all 8 essential handler methods: `handle_start_tag`, `handle_end_tag`, `handle_data`, `handle_entity_ref`, `handle_entity_decl`, `handle_comment`, `handle_pi`, and `handle_decl`.

### Example 8.8. Introducing `HTMLParser`

```
class SimpleHTMLParser(HTMLParser):
 """A simple HTML parser that prints out the tags and data it finds.
 It is a subclass of HTMLParser and implements the 8 essential
 handler methods. It prints out the tag name and the attributes for
 start tags, the tag name for end tags, the data for data tags, the
 entity reference for entity references, the entity declaration for
 entity declarations, the comment for comments, the processing
 instruction for processing instructions, and the declaration for
 declarations. It also prints out the text data it finds between
 tags. The parser is case-insensitive and does not check for
 well-formedness. It is a simple parser and should not be used
 for anything more than educational purposes.
 """
 def __init__(self):
 super().__init__()
 self.reset()
 self.tag_stack = []
 self.data = []
 self.entity_ref = []
 self.entity_decl = []
 self.comment = []
 self.pi = []
 self.decl = []

 def handle_start_tag(self, tag, attrs):
 print(f'Start tag: {tag} {attrs}')
 self.tag_stack.append(tag)

 def handle_end_tag(self, tag):
 print(f'End tag: {tag}')
 self.tag_stack.pop()

 def handle_data(self, data):
 print(f'Data: {data}')
 self.data.append(data)

 def handle_entity_ref(self, data):
 print(f'Entity reference: {data}')
 self.entity_ref.append(data)

 def handle_entity_decl(self, data):
 print(f'Entity declaration: {data}')
 self.entity_decl.append(data)

 def handle_comment(self, data):
 print(f'Comment: {data}')
 self.comment.append(data)

 def handle_pi(self, data):
 print(f'Processing instruction: {data}')
 self.pi.append(data)

 def handle_decl(self, data):
 print(f'Declaration: {data}')
 self.decl.append(data)
```

- 1 Called by `HTMLParser.__init__` initializes `self.tag_stack` as an empty list before calling the ancestor method. `self.data` is a data attribute which will hold the pieces of the HTML document.
- 2 Since `SimpleHTMLParser` does not define any methods for specific tags (like the `handle_start_tag` method in `HTMLParser`), it will call `HTMLParser.handle_start_tag` for every start tag. This method takes the tag (`tag`) and the attributes (`attrs`) and prints them out. It also appends the tag name to `self.tag_stack`.
- 3 Reconstructing end tags is much simpler; just take the tag name and wrap it in the `</>` brackets.
- 4 When `SimpleHTMLParser` finds a character reference, it calls `HTMLParser.handle_entity_ref` with the bare reference. If the HTML document contains the reference `&#x26;` it will be `&`. Reconstructing the original character reference requires wrapping the reference in `&#x26;` characters.
- 5 Entity references are similar to character references, but without the hash mark. Reconstructing the original entity reference requires wrapping the reference in `&` characters.
- 6 Blocks of text are simply appended to `self.data` unaltered.
- 7 HTML comments are wrapped in `<!-->` characters.
- 8 Processing instructions are wrapped in `<?>` characters.

The HTML specification requires  that all non-HTML (like client-side JavaScript) must be enclosed in HTML comments, but not all web pages follow this rule.

## Example 8.9. Output



- 1 This is the one method in `HTMLParser` that is never called by the ancestor `BaseHTMLParser`. Since the other handler methods store their reconstructed HTML in `self.asf`, this method is the only one that can be overridden.
- 2 If you prefer, you could use the `unescape` method of the `HTMLParser` module instead: `parser.unescape(html)`.

### Further reading

W3C (<http://www.w3.org/>) discusses character and entity references (<http://www.w3.org/TR/REC-html40/charset.html#entities>).  
*Python Library Reference* (<http://www.python.org/doc/current/lib/>) confirms your suspicions that the `HTMLParser` module (<http://www.python.org/doc/current/lib/module-htmlparser.html>) is the one to use.

## 8.5. and

Let's digress from HTML processing for a minute and talk about how Python handles variables. Python has two built-in functions, `and` and `or`, which are used to combine boolean expressions.

Remember `and`? You first saw it here:



No, wait, you can't learn about `and` yet. First, you need to learn about namespaces. This is dry stuff, but it's important, so pay attention.

Python uses what are called namespaces to keep track of variables. A namespace is just like a dictionary where the keys are names of variables and the values are the objects they refer to.

At any particular point in a Python program, there are several namespaces available. Each function has its own namespace, called the local namespace.

When a line of code asks for the value of a variable `x`, Python will search for that variable in all the available namespaces, in order:

- local namespace - specific to the current function or class method. If the function defines a local variable `x` or has an argument `x`, Python will use that value.
- global namespace - specific to the current module. If the module has defined a variable, function, or class called `x`, Python will use that value.
- built-in namespace - global to all modules. As a last resort, Python will assume that `x` is the name of built-in function or variable.

If Python doesn't find `x` in any of these namespaces, it gives up and raises a `NameError` with the message `name 'x' is not defined`, which you saw back in Example 3.18, "Referencing a Variable That Doesn't Exist".

Python 2.2 introduced a subtle but important change that affects the namespace search order: nested scopes. In versions of Python prior to 2.2, the search order was global, local, and built-in. In Python 2.2 and later, the search order is local, global, and built-in.

Are you confused yet? Don't despair! This is really cool, I promise. Like many things in Python, namespaces are *directly accessible at run-time*.

## Example 8.10. Introducing



- 1 The function `vars` has two variables in its local namespace: `g` whose value is passed in to the function, and `x` which is defined within the function.
- 2 `vars` returns a dictionary of name/value pairs. The keys of this dictionary are the names of the variables as strings; the values of the dictionary are the values of the variables.
- 3 Remember, Python has dynamic typing, so you could just as easily pass a string in for `g` the function (and the call to `vars` would still work).

What does `vars` do for the module namespace? It does the same thing for the global namespace. It returns all the functions and classes defined in the module.

Remember the difference between `import` and `from`. With `import` the module itself is imported, but it retains its own namespace, which is why you need to use `module.name`.

## Example 8.11. Introducing `vars`

Look at the following block of code at the bottom of `vars.py`:

```
def vars(g):
 return {'g': g, 'x': x}
```

- 1 Just so you don't get intimidated, remember that you've seen all this before. The `vars` function returns a dictionary, and you're iterating through it.

Now running the script from the command line gives this output (note that your output may be slightly different, depending on your platform):

```
>>> from vars import vars
>>> vars(vars)
{'g': <module 'vars' from 'vars.py'>, 'x': 1}
```

- 1 `vars` was imported from `vars` using `from`. That means that it was imported directly into the module's namespace, and here it is.
- 2 Contrast this with `import vars` which was imported using `import`. That means that the `vars` module itself is in the namespace, but the `vars` variable defined within the module is not.
- 3 This module only defines one class, `vars`, and here it is. Note that the value here is the class itself, not a specific instance of the class.
- 4 Remember the `__name__` trick? When running a module (as opposed to importing it from another module), the built-in `__name__` attribute is a special value.

Using the `vars` and `globals` functions, you can get the value of arbitrary variables dynamically, providing the variable name as a string. This mirrors the `vars` function.

There is one other important difference between the `vars` and `globals` functions, which you should learn now before it bites you. It will bite you anyway.

## Example 8.12. `vars` is read-only, `globals` is not

```
def vars(g):
 return {'g': g, 'x': x}

def globals():
 return {'g': g, 'x': x}
```

- 1 Since `z` is called with `3` this will print `3`. This should not be a surprise.
- 2 `z` is a function that returns a dictionary, and here you are setting a value in that dictionary. You might think that this would change the value of `z`, but it doesn't.
- 3 This prints `3`, not `2`.
- 4 After being burned by `z` you might think that this *wouldn't* change the value of `z` but it does. Due to internal differences in how Python handles dictionaries, this does change the value of `z`.
- 5 This prints `2`, not `3`.

## 8.6. Dictionary-based string formatting

Why did you learn about `dict` and `z`? So you can learn about dictionary-based string formatting. As you recall, regular string formatting provides a way to format strings using a tuple of values.

There is an alternative form of string formatting that uses dictionaries instead of tuples of values.

### Example 8.13. Introducing dictionary-based string formatting

```
def z():
 return {'x': 3, 'y': 2}
```

- 1 Instead of a tuple of explicit values, this form of string formatting uses a dictionary, `z`. And instead of a simple `%` marker in the string, the format string uses `{x}` to refer to the value of `x` in the dictionary.
- 2 Dictionary-based string formatting works with any number of named keys. Each key must exist in the given dictionary, or the formatting will fail.
- 3 You can even specify the same key twice; each occurrence will be replaced with the same value.

So why would you use dictionary-based string formatting? Well, it does seem like overkill to set up a dictionary of keys and values simply to format a string.

### Example 8.14. Dictionary-based string formatting in `print`

```
print('{x} {y}'.format(x=3, y=2))
```

- 1 Using the built-in `print` function is the most common use of dictionary-based string formatting. It means that you can use the names of local variables directly in the format string.

### Example 8.15. More dictionary-based string formatting

```
def f():
 return ['x', 'y']
```

- 1 When this method is called, `z` is a list of key/value tuples, just like the `dict` of a dictionary, which means you can use multi-variable assignment to unpack the list. Suppose `z` is `[('x', 3), ('y', 2)]`. In the first round of the list comprehension, `x` will get `3` and `y` will get `2`. The string formatting `'{x} {y}'` will resolve to `'3 2'`. This string becomes the first element of the list comprehension's return value. In the second round, `x` will get `2` and `y` will get `3`. The string formatting will resolve to `'2 3'`. The list comprehension returns a list of these two resolved strings, and `z` will join both elements of this list together to form `'3 2 2 3'`.

② Now, using dictionary-based string formatting, you insert the value of `url` into a string. So if `url` is `'http://www.google.com'`, the final result would be `'http://www.google.com'` and that is

Using dictionary-based string formatting with `%` is a convenient way of making complex string formatting expressions more readable, but it can be a bit tedious.

## 8.7. Quoting attribute values

This is a really nice project in Python (<http://groups.google.com/group/html5-a-stamp-of-python>) using a large project of HTML5 cloning.

It consumes HTML (since it's descended from `HTML`) and produces equivalent HTML, but the HTML output is not identical to the input. Tags and attributes are properly quoted.

### Example 8.16. Quoting attribute values

```
1 from __future__ import unicode_literals
2 import sys
3 import re
4 import os
5 import sys
6 import re
7 import os
8 import sys
9 import re
10 import os
11 import sys
12 import re
13 import os
14 import sys
15 import re
16 import os
17 import sys
18 import re
19 import os
20 import sys
21 import re
22 import os
23 import sys
24 import re
25 import os
26 import sys
27 import re
28 import os
29 import sys
30 import re
31 import os
32 import sys
33 import re
34 import os
35 import sys
36 import re
37 import os
38 import sys
39 import re
40 import os
41 import sys
42 import re
43 import os
44 import sys
45 import re
46 import os
47 import sys
48 import re
49 import os
50 import sys
51 import re
52 import os
53 import sys
54 import re
55 import os
56 import sys
57 import re
58 import os
59 import sys
60 import re
61 import os
62 import sys
63 import re
64 import os
65 import sys
66 import re
67 import os
68 import sys
69 import re
70 import os
71 import sys
72 import re
73 import os
74 import sys
75 import re
76 import os
77 import sys
78 import re
79 import os
80 import sys
81 import re
82 import os
83 import sys
84 import re
85 import os
86 import sys
87 import re
88 import os
89 import sys
90 import re
91 import os
92 import sys
93 import re
94 import os
95 import sys
96 import re
97 import os
98 import sys
99 import re
100 import os
```

- ① Note that the attribute values of the `img` attributes in the `img` tags are not properly quoted. (Also note that you're using triple quotes for some of the strings.)
- ② Feed the parser.
- ③ Using the `render` function defined in `html5lib` you get the output as a single string, complete with quoted attribute values. While this may seem anti-climactic, it's the only way to get the output as a single string.

## 8.8. Introducing `re`

`re` is a simple (and silly) descendant of `re`. It runs blocks of text through a series of substitutions, but it makes sure that anything within a `re` block is properly quoted.

To handle the `re` blocks, you define two methods in `re` and `re`.

### Example 8.17. Handling specific tags

- ```
1 def handle_tag(tag, attributes, html):
2     """Call the appropriate handler for the tag and its attributes,
3     passing the tag name, attribute list, and the tag's contents
4     as arguments. If the tag is an end tag, the contents will be None.
5     """
6     # First, we call the default handler for the tag.
7     # This handler will call the appropriate handler for the tag's
8     # contents, if any.
9     # If the tag is an end tag, the contents will be None.
10    # If the tag is a start tag, the contents will be a list of
11    # elements.
12    # If the tag is a self-closing tag, the contents will be None.
```
- 1 is called every time finds a tag in the HTML source. (In a minute, you'll see exactly how this happens.) The method takes a single parameter, tag, which is the tag name.
 - 2 In the method, you initialize a data attribute that serves as a counter for tags. Every time you hit a tag, you increment the counter; eventually, you'll see how this works.
 - 3 That's it, that's the only special processing you do for tags. Now you pass the list of attributes along to so it can do the default processing.
 - 4 is called every time finds a tag. Since end tags can not contain attributes, the method takes no parameters.
 - 5 First, you want to do the default processing, just like any other end tag.
 - 6 Second, you decrement your counter to signal that this block has been closed.

At this point, it's worth digging a little further into what we've claimed repeatedly (and you've taken it on faith so far) that BeautifulSoup looks for and calls specific methods.

Example 8.18. BeautifulSoup's find method

- ```
1 def find(self, name=None, attrs={}, recursive=True, text_only=False):
2 """Find the first node matching the criteria.
3 """
4 # If the criteria is empty, return the first node.
5 if name is None and attrs == {}:
6 return self.firstChild()
7 # If the criteria is not empty, find the first node matching the criteria.
8 # If the criteria is not empty, find the first node matching the criteria.
9 # If the criteria is not empty, find the first node matching the criteria.
10 # If the criteria is not empty, find the first node matching the criteria.
11 # If the criteria is not empty, find the first node matching the criteria.
12 # If the criteria is not empty, find the first node matching the criteria.
13 # If the criteria is not empty, find the first node matching the criteria.
14 # If the criteria is not empty, find the first node matching the criteria.
15 # If the criteria is not empty, find the first node matching the criteria.
16 # If the criteria is not empty, find the first node matching the criteria.
17 # If the criteria is not empty, find the first node matching the criteria.
18 # If the criteria is not empty, find the first node matching the criteria.
19 # If the criteria is not empty, find the first node matching the criteria.
20 # If the criteria is not empty, find the first node matching the criteria.
21 # If the criteria is not empty, find the first node matching the criteria.
22 # If the criteria is not empty, find the first node matching the criteria.
23 # If the criteria is not empty, find the first node matching the criteria.
24 # If the criteria is not empty, find the first node matching the criteria.
25 # If the criteria is not empty, find the first node matching the criteria.
26 # If the criteria is not empty, find the first node matching the criteria.
27 # If the criteria is not empty, find the first node matching the criteria.
28 # If the criteria is not empty, find the first node matching the criteria.
29 # If the criteria is not empty, find the first node matching the criteria.
30 # If the criteria is not empty, find the first node matching the criteria.
31 # If the criteria is not empty, find the first node matching the criteria.
32 # If the criteria is not empty, find the first node matching the criteria.
33 # If the criteria is not empty, find the first node matching the criteria.
34 # If the criteria is not empty, find the first node matching the criteria.
35 # If the criteria is not empty, find the first node matching the criteria.
36 # If the criteria is not empty, find the first node matching the criteria.
37 # If the criteria is not empty, find the first node matching the criteria.
38 # If the criteria is not empty, find the first node matching the criteria.
39 # If the criteria is not empty, find the first node matching the criteria.
40 # If the criteria is not empty, find the first node matching the criteria.
41 # If the criteria is not empty, find the first node matching the criteria.
42 # If the criteria is not empty, find the first node matching the criteria.
43 # If the criteria is not empty, find the first node matching the criteria.
44 # If the criteria is not empty, find the first node matching the criteria.
45 # If the criteria is not empty, find the first node matching the criteria.
46 # If the criteria is not empty, find the first node matching the criteria.
47 # If the criteria is not empty, find the first node matching the criteria.
48 # If the criteria is not empty, find the first node matching the criteria.
49 # If the criteria is not empty, find the first node matching the criteria.
50 # If the criteria is not empty, find the first node matching the criteria.
51 # If the criteria is not empty, find the first node matching the criteria.
52 # If the criteria is not empty, find the first node matching the criteria.
53 # If the criteria is not empty, find the first node matching the criteria.
54 # If the criteria is not empty, find the first node matching the criteria.
55 # If the criteria is not empty, find the first node matching the criteria.
56 # If the criteria is not empty, find the first node matching the criteria.
57 # If the criteria is not empty, find the first node matching the criteria.
58 # If the criteria is not empty, find the first node matching the criteria.
59 # If the criteria is not empty, find the first node matching the criteria.
60 # If the criteria is not empty, find the first node matching the criteria.
61 # If the criteria is not empty, find the first node matching the criteria.
62 # If the criteria is not empty, find the first node matching the criteria.
63 # If the criteria is not empty, find the first node matching the criteria.
64 # If the criteria is not empty, find the first node matching the criteria.
65 # If the criteria is not empty, find the first node matching the criteria.
66 # If the criteria is not empty, find the first node matching the criteria.
67 # If the criteria is not empty, find the first node matching the criteria.
68 # If the criteria is not empty, find the first node matching the criteria.
69 # If the criteria is not empty, find the first node matching the criteria.
70 # If the criteria is not empty, find the first node matching the criteria.
71 # If the criteria is not empty, find the first node matching the criteria.
72 # If the criteria is not empty, find the first node matching the criteria.
73 # If the criteria is not empty, find the first node matching the criteria.
74 # If the criteria is not empty, find the first node matching the criteria.
75 # If the criteria is not empty, find the first node matching the criteria.
76 # If the criteria is not empty, find the first node matching the criteria.
77 # If the criteria is not empty, find the first node matching the criteria.
78 # If the criteria is not empty, find the first node matching the criteria.
79 # If the criteria is not empty, find the first node matching the criteria.
80 # If the criteria is not empty, find the first node matching the criteria.
81 # If the criteria is not empty, find the first node matching the criteria.
82 # If the criteria is not empty, find the first node matching the criteria.
83 # If the criteria is not empty, find the first node matching the criteria.
84 # If the criteria is not empty, find the first node matching the criteria.
85 # If the criteria is not empty, find the first node matching the criteria.
86 # If the criteria is not empty, find the first node matching the criteria.
87 # If the criteria is not empty, find the first node matching the criteria.
88 # If the criteria is not empty, find the first node matching the criteria.
89 # If the criteria is not empty, find the first node matching the criteria.
90 # If the criteria is not empty, find the first node matching the criteria.
91 # If the criteria is not empty, find the first node matching the criteria.
92 # If the criteria is not empty, find the first node matching the criteria.
93 # If the criteria is not empty, find the first node matching the criteria.
94 # If the criteria is not empty, find the first node matching the criteria.
95 # If the criteria is not empty, find the first node matching the criteria.
96 # If the criteria is not empty, find the first node matching the criteria.
97 # If the criteria is not empty, find the first node matching the criteria.
98 # If the criteria is not empty, find the first node matching the criteria.
99 # If the criteria is not empty, find the first node matching the criteria.
100 # If the criteria is not empty, find the first node matching the criteria.
```
- 1 At this point, BeautifulSoup has already found a start tag and parsed the attribute list. The only thing left to do is figure out whether there is a specific method to call.
  - 2 The "magic" of BeautifulSoup is nothing more than your old friend, getattr. What you may not have realized before is that BeautifulSoup will find methods defined in descendant classes.
  - 3 BeautifulSoup raises an AttributeError if the method it's looking for doesn't exist in the object (or any of its descendants), but that's okay, because you wrapped the call in a try/except block.
  - 4 Since you didn't find a class method, you'll also look for a static method before giving up. This alternate naming scheme is generally used for static methods.
  - 5 Another way to look at this is that the call to getattr failed with AttributeError. Since you found neither a class method nor a static method for this tag, you catch the exception and return None.
  - 6 Remember, BeautifulSoup blocks can have an if clause, which is called if no exception is raised during the block. Logically, that means that you *did* find a method.
  - 7 By the way, don't worry about these different return values; in theory they mean something, but they're never actually used. Don't worry about them.
  - 8 BeautifulSoup and BeautifulSoup methods are not called directly; the tag, method, and attributes are passed to this function, BeautifulSoup.find, so that descendants can override it and call their own find method.

Now back to our regularly scheduled program: When you left, you were in the process of defining specific handler methods for and tags. To

### Example 8.19. Overriding the method



- 1 is called with only one argument, the text to process.
- 2 In the ancestor the method simply appended the text to the output buffer, Here the logic is only slightly more complicated. If you're

You're close to completely understanding The only missing link is the nature of the text substitutions themselves. If you know any Perl, you

## 8.9. Putting it all together

It's time to put everything you've learned so far to good use. I hope you were paying attention.

### Example 8.20. The function, part 1



- 1 The function has an optional argument which is a string that specifies the dialect you'll be using. You'll see how this is used in a minute.
- 2 Hey, wait a minute, there's an statement in this function! That's perfectly legal in Python. You're used to seeing statements at the top
- 3 Now you get the source of the given URL.

### Example 8.21. The function, part 2: curiouuser and curiouuser



- 1 is a string method you haven't seen before; it simply capitalizes the first letter of a string and forces everything else to lowercase. Con
- 2 You have the name of a class as a string (c) and you have the global namespace as a dictionary (g). Combined, you can get a reference
- 3 Finally, you have a class object (C) and you want an instance of the class. Well, you already know how to do that: call the class like a f

Why bother? After all, there are only 3 classes; why not just use a statement? (Well, there's no statement in Python, but why not just use a s

Even better, imagine putting in a separate module, and importing it with . You've already seen that this includes it in g, so b would still work

Now imagine that the name of the dialect is coming from somewhere outside the program, maybe from a database or from a user-inputted va

Finally, imagine a framework with a plug-in architecture. You could put each class in a separate file, leaving only the function in . Assuming

### Example 8.22. The function, part 3



- ❶ After all that imagining, this is going to seem pretty boring, but the `consume` function is what does the entire transformation. You had the entire
- ❷ Because `consume` maintains an internal buffer, you should always call the parser's `close` method when you're done (even if you fed it all at once, like
- ❸ Remember, `consume` is the function you defined on `Parser` that joins all the pieces of output you've buffered and returns them in a single string.

And just like that, you've "translated" a web page, given nothing but a URL and the name of a dialect.

## Further reading

You thought I was kidding about the server-side scripting idea. So did I, until I found this web-based dialectizer (<http://rinkworks.com>).

## 8.10. Summary

Python provides you with a powerful tool, `BeautifulSoup`, to manipulate HTML by turning its structure into an object model. You can use this tool in many

- parsing the HTML looking for something specific·

- aggregating the results, like the URL lister·

- altering the structure along the way, like the attribute quoter·

- transforming the HTML into something else by manipulating the text while leaving the tags alone, like the `HTML2Text`·

Along with these examples, you should be comfortable doing all of the following things:

- Using `from .. import *` and `import .. as ..` to access namespaces·

- Formatting strings using dictionary-based substitutions·

---

[T]he technical term for a parser like `BeautifulSoup` is a *consumer*: it consumes HTML and breaks it down. Presumably, the name `consume` was chosen to fit into the

[T]he reason Python is better at lists than strings is that lists are mutable but strings are immutable. This means that appending to a list just adds

[B]ut don't get out much.

[A]ll right, it's not that common a question. It's not up there with "What editor should I use to write Python code?" (answer: Emacs) or "Is Python



## Chapter 9. XML Processing

### 9.1. Diving in

These next two chapters are about XML processing in Python. It would be helpful if you already knew what an XML document looks like, th

If you're not particularly interested in XML, you should still read these chapters, which cover important topics like Python packages, Unico

Being a philosophy major is not required, although if you have ever had the misfortune of being subjected to the writings of Immanuel Kant,

There are two basic ways to work with XML. One is called SAX ("Simple API for XML"), and it works by reading the XML a little bit at a t

The following is a complete Python program which generates pseudo-random output based on a context-free grammar defined in an XML fo

#### Example 9.1.

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5>).





















































```
def __init__(self, data):
 self.data = data

 def __str__(self):
 return str(self.data)

 def __repr__(self):
 return repr(self.data)

 def __len__(self):
 return len(self.data)

 def __getitem__(self, index):
 return self.data[index]

 def __setitem__(self, index, value):
 self.data[index] = value

 def __delitem__(self, index):
 del self.data[index]

 def __iter__(self):
 return iter(self.data)

 def __contains__(self, value):
 return value in self.data

 def __add__(self, other):
 return self.data + other

 def __sub__(self, other):
 return self.data - other

 def __mul__(self, other):
 return self.data * other

 def __div__(self, other):
 return self.data / other

 def __mod__(self, other):
 return self.data % other

 def __pow__(self, other):
 return self.data ** other

 def __lt__(self, other):
 return self.data < other

 def __le__(self, other):
 return self.data <= other

 def __gt__(self, other):
 return self.data > other

 def __ge__(self, other):
 return self.data >= other

 def __eq__(self, other):
 return self.data == other

 def __ne__(self, other):
 return self.data != other

 def __hash__(self):
 return hash(self.data)
```

一、二、三、四、五、六、七、八、九、十、十一、十二、十三、十四、十五、十六、十七、十八、十九、二十、二十一、二十二、二十三、二十四、二十五、二十六、二十七、二十八、二十九、三十、三十一、三十二、三十三、三十四、三十五、三十六、三十七、三十八、三十九、四十、四十一、四十二、四十三、四十四、四十五、四十六、四十七、四十八、四十九、五十、五十一、五十二、五十三、五十四、五十五、五十六、五十七、五十八、五十九、六十、六十一、六十二、六十三、六十四、六十五、六十六、六十七、六十八、六十九、七十、七十一、七十二、七十三、七十四、七十五、七十六、七十七、七十八、七十九、八十、八十一、八十二、八十三、八十四、八十五、八十六、八十七、八十八、八十九、九十、九十一、九十二、九十三、九十四、九十五、九十六、九十七、九十八、九十九、一百。



[illegible]

### Example 9.2.







Let me repeat that this is much, much funnier if you are now or have ever been a philosophy major.

### Example 9.4. Simpler output from

You will take a closer look at the structure of the grammar file later in this chapter. For now, all you need to know is that the grammar file de

## 9.2. Packages

### Example 9.5. Loading an XML document (a sneak peek)

# Dive Into Python





If you ever find yourself writing a large subsystem in Python (or, more likely, when you realize that your small subsystem has grown into a large one), you may want to consider using XML to represent your data.

## 9.3. Parsing XML

As I was saying, actually parsing an XML document is very simple: one line of code. Where you go from there is up to you.

### Example 9.8. Loading an XML document (for real this time)

```
1 from xml.etree import ElementTree
2
3 # Load the XML document
4 tree = ElementTree.parse('example.xml')
5
6 # Print out the XML
7 print(tree.getroot().getchildren())
```

1 As you saw in the previous section, this imports the `ElementTree` module from the `xml.etree` package.

2 Here is the one line of code that does all the work: `parse` takes one argument and returns a parsed representation of the XML document. The

3 The object returned from `parse` is a `Tree` object, a descendant of the `Element` class. This `Tree` object is the root level of a complex tree-like structure of interlocking

4 `Element` is a method of the `Tree` class (and is therefore available on the `Tree` object you got from `parse`) prints out the XML that this `Tree` represents. For the `Tree`

Now that you have an XML document in memory, you can start traversing through it.

### Example 9.9. Getting child nodes

```
1 # Get the first child node
2 root = tree.getroot()
3 first_child = root.getchildren()[0]
```

1 Every `Element` has a `children` attribute, which is a list of the `Element` objects. A `Tree` always has only one child node, the root element of the XML document (in this case, the `root` element).

2 To get the first (and in this case, the only) child node, just use regular list syntax. Remember, there is nothing special going on here; the `children` attribute is just a list.

3 Since getting the first child node of a node is a useful and common activity, the `Element` class has a `getchildren` attribute, which is synonymous with `children`. (The `getchildren` attribute is a method, not a property, so you need to use parentheses to call it.)

### Example 9.10. `getchildren` works on any node

```
1 # Get the first child node
2 root = tree.getroot()
3 first_child = root.getchildren()[0]
```



- ① Since the `find` method is defined in the `Element` class, it is available on any XML node, not just the `Element`.

**Example 9.11. Child nodes can be text**



- ① Looking at the XML in `example9_11.xml` you might think that the `root` has only two child nodes, the two `Element`s. But you're missing something: the carriage return after the `<?xml` tag and before the first `<tag`.
- ② The first child is a `Text` object representing the carriage return after the `<?xml` tag and before the first `<tag`.
- ③ The second child is an `Element` object representing the first `Element`.
- ④ The fourth child is an `Element` object representing the second `Element`.
- ⑤ The last child is a `Text` object representing the carriage return after the `</tag>` tag and before the `</tag>` tag.

**Example 9.12. Drilling down all the way to text**





- 1 As you saw in the previous example, the first `Element` is `None` since `childNodes[0]` is a `Text` node for the carriage return.
- 2 The `Element` has its own set of child nodes, one for the carriage return, a separate one for the spaces, one for the `Element`, and so forth.
- 3 You can even use the `find` method here, deeply nested within the document.
- 4 The `Element` has only one child node (you can't tell that from this example, but look at `childNodes` if you don't believe me), and it is a `Text` node for the text.
- 5 The `data` attribute of a `Text` node gives you the actual string that the text node represents. But what is that `un` front of the string? The answer to

## 9.4. Unicode

Unicode is a system to represent characters from all the world's different languages. When Python parses an XML document, all data is stored in Unicode. You'll get to all that in a minute, but first, some background.

**Historical note.** Before unicode, there were separate character encoding systems for each language, each using the same numbers (0-255) to represent different characters. This caused a problem: how to represent a character that is not in the original set of characters. The solution was to use a 2-byte number, from 0 to 65535, to represent a character that is not in the original set of characters.

Of course, there is still the matter of all these legacy encoding systems. 7-bit ASCII, for instance, which stores English characters as numbers 0-127. When dealing with unicode data, you may at some point need to convert the data back into one of these other legacy encoding systems. For instance, to store it in a database. And on that note, let's get back to Python.

Python has had unicode support throughout the language since version 2.0. The XML package uses unicode to store all parsed XML data, but

### Example 9.13. Introducing unicode



- 1 To create a unicode string instead of a regular ASCII string, add the letter "u" before the string. Note that this particular string doesn't have any non-ASCII characters.
- 2 When printing a string, Python will attempt to convert it to your default encoding, which is usually ASCII. (More on this in a minute.)

### Example 9.14. Storing non-ASCII characters







- ❶ This is a sample extract from a real Russian XML document; it's part of a Russian translation of this very book. Note the encoding, `utf-8`.
- ❷ These are Cyrillic characters which, as far as I know, spell the Russian word for "Preface". If you open this file in a regular text editor,

### Example 9.19. Parsing



Предисловие

- ❶ I'm assuming here that you saved the previous example as `example9_19.py` in the current directory. I am also, for the sake of completeness, assuming
- ❷ Note that the text data of the `title` tag (now in the `title` variable, thanks to that long concatenation of Python functions which I hastily skipped o
- ❸ Printing the title is not possible, because this unicode string contains non-ASCII characters, so Python can't convert it to ASCII becau
- ❹ You can, however, explicitly convert it to `str` in which case you get a (regular, not unicode) string of single-byte characters (`025` and so f
- ❺ Printing the `str` encoded string will probably show gibberish on your screen, because your Python IDE is interpreting those characters as

To sum up, unicode itself is a bit intimidating if you've never seen it before, but unicode data is really very easy to handle in Python. If your

### Further reading

Unicode.org (<http://www.unicode.org/>) is the home page of the unicode standard, including a brief technical introduction (<http://www.unicode.org/unicode/5.0/ucd/ucd.html>) and a Unicode Tutorial ([http://www.reportlab.com/i18n/python\\_unicode\\_tutorial.html](http://www.reportlab.com/i18n/python_unicode_tutorial.html)) has some more examples of how to use Python's u

PEP 263 (<http://www.python.org/peps/pep-0263.html>) goes into more detail about how and when to define a character encoding in y

## 9.5. Searching for elements

Traversing XML documents by stepping through each node can be tedious. If you're looking for something in particular, buried deep within

For this section, you'll be using the `grammar` file, which looks like this:

### Example 9.20.



```

<?xml version="1.0" encoding="UTF-8" ?>
<root>
 <bit>0</bit>
 <bit>1</bit>
</root>

```

It has two `bit` and `bit` is either a `0` or `1` and a `8` is

### Example 9.21. Introducing `find`

```

from xml.etree.ElementTree import ElementTree, Element
tree = ElementTree.parse('bits.xml')
root = tree.getroot()
bits = root.findall('bit')

```

- 1 `find` takes one argument, the name of the element you wish to find. It returns a list of `Element` objects, corresponding to the XML elements that have the given name.

### Example 9.22. Every element is searchable

```

from xml.etree.ElementTree import ElementTree, Element
tree = ElementTree.parse('bits.xml')
root = tree.getroot()
bits = root.findall('bit')
for bit in bits:
 print(bit.get('value'))

```

- 1 Continuing from the previous example, the first object in your `bits` list is the `Element`.
- 2 You can use the same `find` method on this `Element` to find all the `elements` within the `Element`.
- 3 Just as before, the `find` method returns a list of all the elements it found. In this case, you have two, one for each bit.

### Example 9.23. Searching is actually recursive

```

from xml.etree.ElementTree import ElementTree, Element
tree = ElementTree.parse('bits.xml')
root = tree.getroot()
bits = root.findall('bit')

```



- 1 Note carefully the difference between this and the previous example. Previously, you were searching for elements within `books`, but here you are searching for elements within `book`.
- 2 The first two elements are within the first `book`.
- 3 The last element is the one within the second `book`.

## 9.6. Accessing element attributes

XML elements can have one or more attributes, and it is incredibly simple to access them once you have parsed an XML document.

For this section, you'll be using the `books.xml` file that you saw in the previous section.

This section may be a little confusing, because of some overlapping terminology. Elements in an XML document have attributes, and Python objects have attributes.

### Example 9.24. Accessing element attributes



- 1 Each `Element` object has an attribute called `attrib` which is a `dict` object. This sounds scary, but it's not, because a `dict` is an object that acts like a dictionary.
- 2 Treating the `attrib` as a dictionary, you can get a list of the names of the attributes of this element by using `attrib.keys()`. This element has only one attribute.
- 3 Attribute names, like all other text in an XML document, are stored in unicode.
- 4 Again treating the `attrib` as a dictionary, you can get a list of the values of the attributes by using `attrib.values()`. The values are themselves objects, of type `unicode`.
- 5 Still treating the `attrib` as a dictionary, you can access an individual attribute by name, using normal dictionary syntax. (Readers who have been following the book since the beginning will recognize this as the same syntax used to access dictionary elements.)

### Example 9.25. Accessing individual attributes



- 1 The `XmlAttribute` object completely represents a single XML attribute of a single XML element. The name of the attribute (the same name as you used in the XML) is stored in `name`.
- 2 The actual text value of this XML attribute is stored in `value`.

Like a dictionary, attributes of an XML element have no ordering. Attributes may *happen to be* listed in a certain order in the original XML document, but that's not guaranteed.

## 9.7. Segue

OK, that's it for the hard-core XML stuff. The next chapter will continue to use these same example programs, but focus on other aspects that are more general to Python.

Before moving on to the next chapter, you should be comfortable doing all of these things:

- Parsing XML documents using `xml.etree.ElementTree`
- Searching through the parsed document, and accessing arbitrary element attributes and element children
- Organizing complex libraries into packages
- Converting unicode strings to different character encodings

---

[This, sadly, is *still* an oversimplification. Unicode now has been extended to handle ancient Chinese, Korean, and Japanese texts, which had



# Chapter 10. Scripts and Streams

## 10.1. Abstracting input sources

One of Python's greatest strengths is its dynamic binding, and one powerful use of dynamic binding is the *file-like object*.

Many functions which require an input source could simply take a filename, go open the file for reading, read it, and close it when they're done.

In the simplest case, a *file-like object* is any object with a `read` method with an optional `size` parameter, which returns a string. When called with no argument, it returns a string of up to 1024 bytes.

This is how reading from real files works; the difference is that you're not limiting yourself to real files. The input source could be anything:

In case you were wondering how this relates to XML processing, `xml.etree.ElementTree.parse` is one such function which can take a file-like object.

### Example 10.1. Parsing XML from a file

```
1 # Parse an XML document from a file
2
3 import sys
4 import xml.etree.ElementTree as ET
5
6 def parse_xml(filename):
7 """Parse an XML document from a file.
8 Returns the root element of the document.
9 """
10 tree = ET.parse(filename)
11 return tree.getroot()
12
13 if __name__ == '__main__':
14 if len(sys.argv) != 2:
15 sys.stderr.write('Usage: %s filename\n' % sys.argv[0])
16 sys.exit(1)
17 root = parse_xml(sys.argv[1])
18 print(root)
```

- 1 First, you open the file on disk. This gives you a file object.
- 2 You pass the file object to `ET.parse` which calls the `read` method of `File` and reads the XML document from the file on disk.
- 3 Be sure to call the `close` method of the file object after you're done with it. `ET.parse` will not do this for you.
- 4 Calling the `getroot` method on the returned XML document prints out the entire thing.

Well, that all seems like a colossal waste of time. After all, you've already seen that `ET.parse` can simply take the filename and do all the opening and closing for you.

### Example 10.2. Parsing XML from a URL

```
1 # Parse an XML document from a URL
2
3 import sys
4 import xml.etree.ElementTree as ET
5
6 def parse_xml(url):
7 """Parse an XML document from a URL.
8 Returns the root element of the document.
9 """
10 tree = ET.parse(url)
11 return tree.getroot()
12
13 if __name__ == '__main__':
14 if len(sys.argv) != 2:
15 sys.stderr.write('Usage: %s url\n' % sys.argv[0])
16 sys.exit(1)
17 root = parse_xml(sys.argv[1])
18 print(root)
```



- 1 As you saw in a previous chapter, `urllib.urlopen` takes a web page URL and returns a file-like object. Most importantly, this object has a `read` method which returns the data from the URL.
- 2 Now you pass the file-like object to `xml.etree.ElementTree.parse` which obediently calls the `read` method of the object and parses the XML data that the `read` method returns.
- 3 As soon as you're done with it, be sure to close the file-like object that `urllib.urlopen` gives you.
- 4 By the way, this URL is real, and it really is XML. It's an XML representation of the current headlines on Slashdot (<http://slashdot.org>).

### Example 10.3. Parsing XML from a string (the easy but inflexible way)



- 1 `xml.etree.ElementTree` has a method, `fromstring`, which takes an entire XML document as a string and parses it. You can use this instead of `parse` if you know you already have a string.
- OK, so you can use the `fromstring` function for parsing both local files and remote URLs, but for parsing strings, you use... a different function. That means you need to use `urllib.urlopen` to get a file-like object, and then pass that object to `parse`. And in fact, there is a module specifically for this purpose.

### Example 10.4. Introducing `xml.etree.ElementTree`



```
from io import StringIO
from xml.parsers.expat import ParserCreate
```

- 1 The `StringIO` module contains a single class, also called `StringIO`, which allows you to turn a string into a file-like object. The `StringIO` class takes the string as a parameter.
- 2 Now you have a file-like object, and you can do all sorts of file-like things with it. Like `read()` which returns the original string.
- 3 Calling `read()` again returns an empty string. This is how real file objects work too; once you read the entire file, you can't read any more without seeking.
- 4 You can explicitly seek to the beginning of the string, just like seeking through a file, by using the `seek()` method of the `StringIO` object.
- 5 You can also read the string in chunks, by passing a `size` parameter to the `read()` method.
- 6 At any time, `read()` will return the rest of the string that you haven't read yet. All of this is exactly how file objects work; hence the term *file-like object*.

### Example 10.5. Parsing XML from a string (the file-like object way)

```
def parse_xml_string(xml_string):
 parser = ParserCreate()
 parser.Parse(xml_string)
```

- 1 Now you can pass the file-like object (really a `StringIO`) to `parse_xml_string()` which will call the object's `read()` method and happily parse away, never knowing that it's a string.

So now you know how to use a single function, `parse_xml_string()`, to parse an XML document stored on a web page, in a local file, or in a hard-coded string. For more on parsing XML, see Chapter 11.

### Example 10.6. `urllib2.urlopen()`

```
import urllib2
url = 'http://www.python.org'
response = urllib2.urlopen(url)
data = response.read()
print data
```

- 1 The `urllib2.urlopen()` function takes a single parameter, `url`, and returns a file-like object. `url` is a string of some sort; it can either be a URL (like `http://www.python.org`) or a path to a file on disk.
- 2 First, you see if `url` is a URL. You do this through brute force: you try to open it as a URL and silently ignore errors caused by trying to open it as a URL.
- 3 On the other hand, if `urllib2.urlopen()` yelled at you and told you that `url` wasn't a valid URL, you assume it's a path to a file on disk and try to open it. Again, this is brute force.

④ By this point, you need to assume that `data` is a string that has hard-coded data in it (since nothing else worked), so you use `print` to create a file

Now you can use this `print_xml` function in conjunction with `open` to make a function that takes a `path` that refers to an XML document somehow (either as a URL

### Example 10.7. Using `print_xml`

```
#!/usr/bin/env python
import sys
import xml.etree.ElementTree as ET

def print_xml(path):
 tree = ET.parse(path)
 root = tree.getroot()
 print ET.tostring(root, encoding='utf-8')
```

## 10.2. Standard input, output, and error

UNIX users are already familiar with the concept of standard input, standard output, and standard error. This section is for the rest of you.

Standard output and standard error (commonly abbreviated `stdout` and `stderr`) are pipes that are built into every UNIX system. When you `print` something, it goes

### Example 10.8. Introducing `stdout` and `stderr`

```
#!/usr/bin/env python
import sys

def main():
 sys.stdout.write('Hello, stdout!\n')
 sys.stderr.write('Hello, stderr!\n')
```

① As you saw in Example 6.9, “Simple Counters”, you can use Python’s built-in `range` function to build simple counter loops that repeat some

② `data` is a file-like object; calling its `write` function will print out whatever string you give it. In fact, this is what the `print` function really does; it adds

③ In the simplest case, `stdout` and `stderr` send their output to the same place: the Python IDE (if you’re in one), or the terminal (if you’re running Python

`stdout` and `stderr` are both file-like objects, like the ones you discussed in Section 10.1, “Abstracting input sources”, but they are both write-only. They have

### Example 10.9. Redirecting output

```
#!/usr/bin/env python
import sys

def main():
 sys.stdout.write('Hello, stdout!\n')
```

(On Windows, you can use `print` instead of `write` to display the contents of a file.)

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5.0/>)

```
#
```



**Q**

③

4

**D**



- 1 This will print to the IDE "Interactive Window" (or the terminal, if running the script from the command line).
- 2 Always save before redirecting it, so you can set it back to normal later.
- 3 Open a file for writing. If the file doesn't exist, it will be created. If the file does exist, it will be overwritten.
- 4 Redirect all further output to the new file you just opened.
- 5 This will be "printed" to the log file only; it will not be visible in the IDE window or on the screen.
- 6 Set back to the way it was before you mucked with it.
- 7 Close the log file.

Redirecting `#` works exactly the same way, using `#` instead of `#`.

### Example 10.10. Redirecting error information



1



If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5>).

#



34

- 1 Open the log file where you want to store debugging information.
- 2 Redirect standard error by assigning the file object of the newly-opened log file to `#`
- 3 Raise an exception. Note from the screen output that this does *not* print anything on screen. All the normal traceback information has been sent to the log file.
- 4 Also note that you're not explicitly closing your log file, nor are you setting `#` back to its original value. This is fine, since once the program exits, the file object is garbage collected and the file is closed.

Since it is so common to write error messages to standard error, there is a shorthand syntax that can be used instead of going through the hassle

### Example 10.11. Printing to `stdout`

6



❶ This is the `find` function from `lxml` which you previously examined in Section 10.1, “Abstracting input sources”. All you’ve done is add three

## 10.3. Caching node lookups

`lxml` employs several tricks which may or may not be useful to you in your XML processing. The first one takes advantage of the consistent structure

A grammar file defines a series of `elements`. Each `element` contains one or more `elements`, which can contain a lot of different things, including `W`

This is how you build up the grammar: define `elements` for the smallest pieces, then define `elements` which “include” the first `elements` by u

This is all very flexible, but there is one downside: performance. When you find an `element` and need to find the corresponding `element`, you have a

### Example 10.14. `lxml`

```
from lxml import etree
import sys

def main():
 # ❶
```

- ❶ Start by creating an empty dictionary, `cache`
- ❷ As you saw in Section 9.5, “Searching for elements”, `find` returns a list of all the elements of a particular name. You easily can get a list of
- ❸ As you saw in Section 9.6, “Accessing element attributes”, you can access individual attributes of an element by name, using standard
- ❹ The values of the `cache` dictionary will be the `elements` themselves. As you saw in Section 9.3, “Parsing XML”, each element, each node, e

Once you build this cache, whenever you come across an `element` and need to find the `element` with the same `attribute`, you can simply look it up in

### Example 10.15. Using the `element` cache

```
def main():
 # ❶
```

You’ll explore the `lxml` function in the next section.

## 10.4. Finding direct children of a node

Another useful technique when parsing XML documents is finding all the direct child elements of a particular element. For instance, in the gra

You might think you could simply use `find` for this, but you can’t. `find` searches recursively and returns a single list for all the elements it finds. Since

### Example 10.16. Finding direct child elements

```
def main():
 # ❶
```

- ❶ As you saw in Example 9.9, “Getting child nodes”, the `getchildren` attribute returns a list of all the child nodes of an element.
- ❷ However, as you saw in Example 9.11, “Child nodes can be text”, the list returned by `getchildren` contains all different types of nodes, including
- ❸ Each node has a `tag` attribute, which can be `None` or any number of other values. The complete list of possible values is in the `lxml` file in the `lxml`

- ④ Once you have a list of actual elements, choosing a random one is easy. Python comes with a module called `random` which includes several useful functions.

## 10.5. Creating separate handlers by node type

The third useful XML processing tip involves separating your code into logical functions, based on node types and element names. Parsed XML documents are represented as Python objects, and each node is a Python object.

### Example 10.17. Class names of parsed XML objects

```
1 from xml.dom.minidom import parse
2
3 # Parse the XML file
4 dom = parse('example.xml')
5
6 # Get the root element
7 root = dom.documentElement
8
9 # Print the class name of the root element
10 print(root.__class__.__name__)
```

- ① Assume for a moment that `example.xml` is in the current directory.
- ② As you saw in Section 9.2, “Packages”, the object returned by parsing an XML document is a `Document` object, as defined in the `xml.dom.minidom` package.
- ③ Furthermore, `__class__` is a built-in attribute of every Python class, and it is a string. This string is not mysterious; it’s the same as the class name.

Fine, so now you can get the class name of any particular XML node (since each XML node is represented as a Python object). How can you use this information?

### Example 10.18. A generic XML node dispatcher

```
1 def dispatch(node):
2 # Dispatch to the appropriate handler
3 handler = getattr(dispatcher, node.__class__.__name__)
4 handler(node)
```

- ① First off, notice that you’re constructing a larger string based on the class name of the node you were passed (in the `node` argument). So if you have a `Text` node, you’ll get `Text`.
- ② Now you can treat that string as a function name, and get a reference to the function itself using `getattr`.
- ③ Finally, you can call that function and pass the node itself as an argument. The next example shows the definitions of each of these functions.

### Example 10.19. Functions called by the dispatcher

```
1 def text_handler(node):
2 # Handle a text node
3 print(node.data)
4
5 def element_handler(node):
6 # Handle an element node
7 print(node.tagName)
8
9 def attribute_handler(node):
10 # Handle an attribute node
11 print(node.name, node.value)
12
13 dispatcher = {
14 'Text': text_handler,
15 'Element': element_handler,
16 'Attribute': attribute_handler,
17 }
```



```
#
```

- 1 `processText` is only ever called once, since there is only one `Node` in an XML document, and only one `Node` object in the parsed XML representation. In this example, `processText` is called on nodes that represent bits of text. The function itself does some special processing to handle automatic capitalization of the text.
- 2 `processText` is just a placeholder since you don't care about embedded comments in the grammar files. Note, however, that you still need to define the function.
- 3 The `processText` method is actually itself a dispatcher, based on the name of the element's tag. The basic idea is the same: take what distinguishes

In this example, the dispatch functions `processText` and `processElement` simply find other methods in the same class. If your processing is very complex (or you have many methods), you might want to use a dispatcher class.

## 10.6. Handling command-line arguments

Python fully supports creating programs that can be run on the command line, complete with command-line arguments and either short- or long-option flags.

It's difficult to talk about command-line processing without understanding how command-line arguments are exposed to your Python program.

### Example 10.20. Introducing `sys.argv`

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5.0.tar.gz>).

```
#!/usr/bin/env python
```

```
import sys
```

```
if __name__ == '__main__':
```

```
 for arg in sys.argv[1:]:
```

- 1 Each command-line argument passed to the program will be in `sys.argv`, which is just a list. Here you are printing each argument on a separate line.

### Example 10.21. The contents of `sys.argv`

```
#!/usr/bin/env python
```

```
import sys
```

```
if __name__ == '__main__':
```

```
 print sys.argv
```

```
 for arg in sys.argv[1:]:
```

```
 print arg
```

```
 print sys.argv[0]
```

```
 print sys.argv[1]
```

```
 print sys.argv[2]
```

```
 print sys.argv[3]
```

```
 print sys.argv[4]
```

```
 print sys.argv[5]
```

```
 print sys.argv[6]
```

```
 print sys.argv[7]
```

```
 print sys.argv[8]
```

```
 print sys.argv[9]
```

```
 print sys.argv[10]
```

```
 print sys.argv[11]
```

```
 print sys.argv[12]
```

```
 print sys.argv[13]
```

```
 print sys.argv[14]
```

```
 print sys.argv[15]
```

```
 print sys.argv[16]
```

```
 print sys.argv[17]
```

```
 print sys.argv[18]
```

```
 print sys.argv[19]
```

```
 print sys.argv[20]
```

```
 print sys.argv[21]
```

- 1 The first thing to know about `sys.argv` is that it contains the name of the script you're calling. You will actually use this knowledge to your advantage.
- 2 Command-line arguments are separated by spaces, and each shows up as a separate element in the `sys.argv` list.
- 3 Command-line flags, like `-h`, also show up as their own element in the `sys.argv` list.
- 4 To make things even more interesting, some command-line flags themselves take arguments. For instance, here you have a flag (`-h`) which

So as you can see, you certainly have all the information passed on the command line, but then again, it doesn't look like it's going to be all that useful.

## Example 10.22. Introducing `argparse`

```
#!/usr/bin/env python
coding: utf-8
'''
 A simple script that demonstrates the use of the argparse module.
 It takes a grammar file or URL as input and prints the usage summary.
'''
import argparse
import sys

def main():
 parser = argparse.ArgumentParser(
 description='A simple script that demonstrates the use of the argparse module. It takes a grammar file or URL as input and prints the usage summary.'
)
 parser.add_argument(
 '-g', '--grammar',
 help='use specified grammar file or URL'
)
 parser.add_argument(
 '-d', '--debug',
 help='show debugging information while parsing'
)
 args = parser.parse_args()

 # Print the usage summary
 parser.print_help()

 # Print the grammar file or URL
 if args.grammar:
 print(args.grammar)

 # Print the debug information
 if args.debug:
 print('Debugging information while parsing')

if __name__ == '__main__':
 main()
```

- 1 First off, look at the bottom of the example and notice that you're calling the `main` function with `if __name__ == '__main__': main()`. Remember, `__name__` is the name of the script that is being executed.
- 2 This is where all the interesting processing happens. The `main` function of the `__main__` module takes three parameters: the argument list (which you can access via `sys.argv`), the parser object, and the `args` object.
- 3 If anything goes wrong trying to parse these command-line flags, `argparse` will raise an exception, which you catch. You told `argparse` all the flags you want to use.
- 4 As is standard practice in the UNIX world, when the script is passed flags it doesn't understand, you print out a summary of proper usage.

So what are all those parameters you pass to the `main` function? Well, the first one is simply the raw list of command-line flags and arguments (not including the script name).

```
h
 print usage summary

g
 use specified grammar file or URL

d
 show debugging information while parsing
```

The first and third flags are simply standalone flags; you specify them or you don't, and they do things (print help) or change state (turn on debug).

To further complicate things, the script accepts either short flags (like `-h`) or long flags (like `--help`) and you want them to do the same thing. This is where `argparse` comes in.

```
h
 print usage summary

g
 use specified grammar file or URL
```

Three things of note here:

- All long flags are preceded by two dashes on the command line, but you don't include those dashes when calling `add_argument`. They are understood by `argparse`.
- The `-g` flag must always be followed by an additional argument, just like the `-g` flag. This is notated by an equals sign, `-g=`.
- The list of long flags is shorter than the list of short flags, because the `-d` flag does not have a corresponding long version. This is fine; it's not required.

Confused yet? Let's look at the actual code and see if it makes sense in context.

## Example 10.23. Handling command-line arguments in `argparse`

```
#!/usr/bin/env python
coding: utf-8
'''
 A simple script that demonstrates the use of the argparse module.
 It takes a grammar file or URL as input and prints the usage summary.
'''
import argparse
import sys

def main():
 parser = argparse.ArgumentParser(
 description='A simple script that demonstrates the use of the argparse module. It takes a grammar file or URL as input and prints the usage summary.'
)
 parser.add_argument(
 '-g', '--grammar',
 help='use specified grammar file or URL'
)
 parser.add_argument(
 '-d', '--debug',
 help='show debugging information while parsing'
)
 args = parser.parse_args()

 # Print the usage summary
 parser.print_help()

 # Print the grammar file or URL
 if args.grammar:
 print(args.grammar)

 # Print the debug information
 if args.debug:
 print('Debugging information while parsing')

if __name__ == '__main__':
 main()
```

ㄚ ㄛ ㄜ ㄝ ㄞ ㄟ ㄠ ㄡ ㄢ ㄣ ㄤ ㄥ ㄦ ㄨ ㄩ ㄣ  
 ㄤ ㄥ ㄦ ㄨ ㄩ ㄣ ㄤ ㄥ ㄦ ㄨ ㄩ ㄣ ㄤ ㄥ ㄦ ㄨ ㄩ ㄣ  
 ㄤ ㄥ ㄦ ㄨ ㄩ ㄣ ㄤ ㄥ ㄦ ㄨ ㄩ ㄣ ㄤ ㄥ ㄦ ㄨ ㄩ ㄣ

- 1 The `g` variable will keep track of the grammar file you're using. You initialize it here in case it's not specified on the command line (using `g = None`).
- 2 The `p` variable that you get back from `p` contains a list of tuples: `g` and `h`. If the flag doesn't take an argument, then `g` will simply be `h`. This means that `g` and `h` are the same object.
- 3 `p` validates that the command-line flags are acceptable, but it doesn't do any sort of conversion between short and long flags. If you specify `-g` and `-h`, `p` will raise an error.
- 4 Remember, the `d` flag didn't have a corresponding long flag, so you only need to check for the short form. If you find it, you set a global `d` to `True`.
- 5 If you find a grammar file, either with a `g` flag or a `h` flag, you save the argument that followed it (stored in `g`) into the `g` variable, overwriting the `g = None` that you set at the top.
- 6 That's it. You've looped through and dealt with all the command-line flags. That means that anything left must be command-line arguments.

## 10.7. Putting it all together

You've covered a lot of ground. Let's step back and see how all the pieces fit together.

To start with, this is a script that takes its arguments on the command line, using the `sys` module.

● ● ● ● ●

You create a new instance of the `Parser` class, and pass it the grammar file and source that may or may not have been specified on the command line:



The `Instance` automatically loads the grammar, which is an XML file. You use your custom `Instance` function to open the file (which could be stored



Oh, and along the way, you take advantage of your knowledge of the structure of the XML document to set up a little cache of references, which

```
def
def
def
```

If you specified some source material on the command line, you use that; otherwise you rip through the grammar looking for the "top-level" node

```
def
def
def
def
def
def
def
```

Now you rip through the source material. The source material is also XML, and you parse it one node at a time. To keep the code separated a

```
def
def
def
```

You bounce through the grammar, parsing all the children of each element,

```
def
.
def
def
```

replacing elements with a random child,

```
def
def
```

and replacing elements with a random child of the corresponding element, which you previously cached.

```
def
def
def
```

Eventually, you parse your way down to plain text,

```
def
def
.
def
```

which you print out.

```
def
.
def
def
```

## 10.8. Summary

Python comes with powerful libraries for parsing and manipulating XML documents. The `lxml` takes an XML file and parses it into Python objects.

Before moving on to the next chapter, you should be comfortable doing all of these things:

- Chaining programs with standard input and output
- Defining dynamic dispatchers with `dispatch`
- Using command-line flags and validating them with `argparse`

# Chapter 11. HTTP Web Services

## 11.1. Diving in

You've learned about HTML processing and XML processing, and along the way you saw how to download a web page and how to parse XML.

Simply stated, HTTP web services are programmatic ways of sending and receiving data from remote servers using the operations of HTTP client and server.

The main advantage of this approach is simplicity, and its simplicity has proven popular with a lot of different sites. Data -- usually XML data -- is sent over HTTP.

Examples of pure XML-over-HTTP web services:

- Amazon API (<http://www.amazon.com/webservices>) allows you to retrieve product information from the Amazon.com online store.

- National Weather Service (<http://www.nws.noaa.gov/alerts/>) (United States) and Hong Kong Observatory (<http://demo.xml.weather.gov.hk/>)

- Atom API (<http://atomenabled.org/>) for managing web-based content.

- Syndicated feeds (<http://syndic8.com/>) from weblogs and news sites bring you up-to-the-minute news from a variety of sites.

In later chapters, you'll explore APIs which use HTTP as a transport for sending and receiving data, but don't map application semantics to the data.

Here is a more advanced version of the `urllib` module that you saw in the previous chapter:

### Example 11.1. `urllib`

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5.0.tar.gz>).

```
#!/usr/bin/env python
import sys
import urllib
import urllib2
import urllib3
import urllib4
import urllib5
import urllib6
import urllib7
import urllib8
import urllib9
import urllib10
import urllib11
import urllib12
import urllib13
import urllib14
import urllib15
import urllib16
import urllib17
import urllib18
import urllib19
import urllib20
import urllib21
import urllib22
import urllib23
import urllib24
import urllib25
import urllib26
import urllib27
import urllib28
import urllib29
import urllib30
import urllib31
import urllib32
import urllib33
import urllib34
import urllib35
import urllib36
import urllib37
import urllib38
import urllib39
import urllib40
import urllib41
import urllib42
import urllib43
import urllib44
import urllib45
import urllib46
import urllib47
import urllib48
import urllib49
import urllib50
import urllib51
import urllib52
import urllib53
import urllib54
import urllib55
import urllib56
import urllib57
import urllib58
import urllib59
import urllib60
import urllib61
import urllib62
import urllib63
import urllib64
import urllib65
import urllib66
import urllib67
import urllib68
import urllib69
import urllib70
import urllib71
import urllib72
import urllib73
import urllib74
import urllib75
import urllib76
import urllib77
import urllib78
import urllib79
import urllib80
import urllib81
import urllib82
import urllib83
import urllib84
import urllib85
import urllib86
import urllib87
import urllib88
import urllib89
import urllib90
import urllib91
import urllib92
import urllib93
import urllib94
import urllib95
import urllib96
import urllib97
import urllib98
import urllib99
import urllib100
```

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99  
100



### Further reading

Paul Prescod believes that pure HTTP web services are the future of the Internet (<http://webservicex.xml.com/pub/a/ws/2002/02/06/>)

## 11.2. How not to fetch data over HTTP

Let's say you want to download a resource over HTTP, such as a syndicated Atom feed. But you don't just want to download it once; you want to download it every time you need it.

### Example 11.2. Downloading a feed the quick-and-dirty way



❶ Downloading anything over HTTP is incredibly easy in Python; in fact, it's a one-liner. The `urllib` module has a handy `urlopen` function that takes a URL and returns a file-like object that you can read from.

So what's wrong with this? Well, for a quick one-off during testing or development, there's nothing wrong with it. I do it all the time. I want to download a feed, and I want to know what's in it. I want to know what's in it.

Let's talk about some of the basic features of HTTP.

## 11.3. Features of HTTP

There are five important features of HTTP which you should support.

### 11.3.1. User-Agent

The `User-Agent` is simply a way for a client to tell a server who it is when it requests a web page, a syndicated feed, or any sort of web service over HTTP.

By default, Python sends a generic `User-Agent`. In the next section, you'll see how to change this to something more specific.



## 11.3.2. Redirects

Sometimes resources move around. Web sites get reorganized, pages move to new addresses. Even web services can reorganize. A syndicated feed can move to a new URL. Every time you request any kind of resource from an HTTP server, the server includes a status code in its response. Status code 301 means "permanent redirect". HTTP has two different ways of signifying that a resource has moved. Status code 302 is a *temporary redirect*; it means "oops, that got moved over there". Your browser will automatically "follow" redirects when it receives the appropriate status code from the HTTP server, but unfortunately, it doesn't tell you where it went.

## 11.3.3. Last-Modified

Some data changes all the time. The home page of CNN.com is constantly updating every few minutes. On the other hand, the home page of a static web site doesn't change. If you ask for the same data a second time (or third, or fourth), you can tell the server the last-modified date that you got last time: you send a `Last-Modified` header. All modern web browsers support last-modified date checking. If you've ever visited a page, re-visited the same page a day later and found that it had changed, that's because the browser checked the last-modified date. Python's URL library has no built-in support for last-modified date checking, but since you can add arbitrary headers to each request and read the response headers, you can do it.

## 11.3.4. ETags

ETags are an alternate way to accomplish the same thing as the last-modified date checking: don't re-download data that hasn't changed. The server sends an `ETag` header with the data, and you send it back in your request. Python's URL library has no built-in support for ETags, but you'll see how to add it later in this chapter.

## 11.3.5. Compression

The last important HTTP feature is gzip compression. When you talk about HTTP web services, you're almost always talking about moving data around. Python's URL library has no built-in support for gzip compression per se, but you can add arbitrary headers to the request. And Python comes with a `gzip` module. Note that our little one-line script to download a syndicated feed did not support any of these HTTP features. Let's see how you can improve it.

## 11.4. Debugging HTTP web services

First, let's turn on the debugging features of Python's HTTP library and see what's being sent over the wire. This will be useful throughout the chapter.

### Example 11.3. Debugging HTTP





- ❶ `urllib` relies on another standard Python library, `urllib`. Normally you don't need to `urllib` directly (`urllib` does that automatically), but you will here so you can see what's going on.
- ❷ Now that the debugging flag is set, information on the the HTTP request and response is printed out in real time. The first thing it tells you is the URL you're requesting.
- ❸ When you request the Atom feed, `urllib` sends three lines to the server. The first line specifies the HTTP verb you're using, and the path of the resource you're requesting.
- ❹ The second line is the `Host` header, which specifies the domain name of the service you're accessing. This is important, because a single HTTP server can serve multiple domain names.
- ❺ The third line is the `User-Agent` header. What you see here is the generic `User-Agent` that the `urllib` library adds by default. In the next section, you'll see how to customize this.
- ❻ The server replies with a status code and a bunch of headers (and possibly some data, which got stored in the `data` variable). The status code is 200, which means "OK".
- ❼ The server tells you when this Atom feed was last modified (in this case, about 13 minutes ago). You can send this date back to the server to avoid requesting the feed again if it hasn't changed.
- ❽ The server also tells you that this Atom feed has an ETag hash of `1234567890`. The hash doesn't mean anything by itself; there's nothing you can do with it.

## 11.5. Setting the `urllib` headers

The first step to improving your HTTP web services client is to identify yourself properly with a `User-Agent`. To do that, you need to move beyond the basic `urllib` client.

### Example 11.4. Introducing `urllib2`



- ❶ If you still have your Python IDE open from the previous section's example, you can skip this, but this turns on HTTP debugging so you can see what's going on.
- ❷ Fetching an HTTP resource with `urllib2` is a three-step process, for good reasons that will become clear shortly. The first step is to create a `urllib2.Request` object.
- ❸ The second step is to build a URL opener. This can take any number of handlers, which control how responses are handled. But you can also use the default opener.

- 4 The final step is to tell the opener to open the URL, using the `obj` you created. As you can see from all the debugging information t

### Example 11.5. Adding headers with the `obj`



- 1 You're continuing from the previous example; you've already created a `obj` with the URL you want to access.
- 2 Using the `add_header` method on the `obj`, you can add arbitrary HTTP headers to the request. The first argument is the header, the second is the value.
- 3 The `obj` you created before can be reused too, and it will retrieve the same feed again, but with your custom header.
- 4 And here's you sending your custom `obj` in place of the generic one that Python sends by default. If you look closely, you'll notice that y

## 11.6. Handling `add_header` and `obj`

Now that you know how to add custom HTTP headers to your web service requests, let's look at adding support for `add_header` and `obj`.

These examples show the output with debugging turned off. If you still have it turned on from the previous section, you can turn it off by sett

### Example 11.6. Testing `add_header`





- ❶ Remember all those HTTP headers you saw printed out when you turned on debugging? This is how you can get access to them programmatically.
- ❷ On the second request, you add the `Last-Modified` header with the last-modified date from the first request. If the data hasn't changed, the server should return a 304 status code.
- ❸ Sure enough, the data hasn't changed. You can see from the traceback that `urllib.error.HTTPError` is raised in response to the 304 status code. `urllib.error.HTTPError` also raises an `urllib.error.HTTPError` exception for conditions that you would think of as errors, such as 404 (page not found). In fact, it will raise `urllib.error.HTTPError` for *any* status code other than 200.

### Example 11.7. Defining URL handlers

This custom URL handler is part of `urllib2`.



- ❶ `urllib2` is designed around URL handlers. Each handler is just a class that can define any number of methods. When something happens -- like a 304 status code -- `urllib2` searches through the defined handlers and calls the `handle_304` method when it encounters a 304 status code from the server. By defining a custom handler, you can intercept the status code and do whatever you want.
- ❷ This is the key part: before returning, you save the status code returned by the HTTP server. This will allow you easy access to it from the `urllib2` object.
- ❸

### Example 11.8. Using custom URL handlers





- 1 You're continuing the previous example, so the `Object` is already set up, and you've already added the `Header`.
- 2 This is the key: now that you've defined your custom URL handler, you need to tell `io` to use it. Remember how I said that `io` broke up the
- 3 Now you can quietly open the resource, and what you get back is an object that, along with the usual headers (use `io` to access them), also
- 4 Note that when the server sends back a `Status` code, it doesn't re-send the data. That's the whole point: to save bandwidth by not re-do

Handling `U` works much the same way, but instead of checking for `U` and sending `U`, you check for `U` and send `U`. Let's start with a fresh IDE session

### Example 11.9. Supporting

- 1 Using the `PseudoDictionary`, you can get the `Entity` returned from the server. (What happens if the server didn't send back an `Entity`? Then this line will throw an exception.)
- 2 OK, you got the data.
- 3 Now set up the second call by setting the `Header` to the `Entity` you got from the first call.
- 4 The second call succeeds quietly (without throwing an exception), and once again you see that the server has sent back a `Status` code.
- 5 Regardless of whether the `Entity`s triggered by `Match` checking or `Hash` matching, you'll never get the data along with the `Entity`. That's the whole point of the `Entity` abstraction.

In these examples, the HTTP server has supported both `Location` and `URI` headers, but not all servers do. As a web services client, you should be prepared

## 11.7. Handling redirects

You can support permanent and temporary redirects using a different kind of custom URL handler.

First, let's see why a redirect handler is necessary in the first place.

### Example 11.10. Accessing web services without a redirect handler





- 1 You'll be better able to see what's happening if you turn on debugging.
- 2 This is a URL which I have set up to permanently redirect to my Atom feed at `http://www.snoke.org/atom.xml`
- 3 Sure enough, when you try to download the data at that address, the server sends back a `301` status code, telling you that the resource has moved.
- 4 The server also sends back a `Location` header that gives the new address of this data.
- 5 `urllib2` notices the redirect status code and automatically tries to retrieve the data at the new location specified in the `Location` header.
- 6 The object you get back from the `urllib2.urlopen` contains the new permanent address and all the headers returned from the second request (retrieved from `http://www.snoke.org/atom.xml`).

This is suboptimal, but easy to fix. `urllib2` doesn't behave exactly as you want it to when it encounters a `301` or `302`. So let's override its behavior. How? Well, we can define a custom redirect handler.

### Example 11.11. Defining the redirect handler

This class is defined in `urllib2.py`:



- 1 Redirect behavior is defined in `urllib2.py` in a class called `RedirectHandler`. You don't want to completely override the behavior, you just want to extend it a little.
- 2 When it encounters a `301` status code from the server, `RedirectHandler` will search through its handlers and call the `handle_301` method. The first thing ours does is just return the status code.
- 3 Here's the key: before you return, you store the status code (in `self.status`) so that the calling program can access it later.
- 4 Temporary redirects (status code `302`) work the same way: override the `handle_302` method, call the ancestor, and save the status code before returning.

So what has this bought us? You can now build a URL opener with the custom redirect handler, and it will still automatically follow redirects.

### Example 11.12. Using the redirect handler to detect permanent redirects






- ❶ First, build a URL opener with the redirect handler you just defined.
- ❷ You sent off a request, and you got a `302` status code in response. At this point, the `redirect` method gets called. You call the ancestor method, which
- ❸ This is the payoff: now, not only do you have access to the new URL, but you have access to the redirect status code, so you can tell the

The same redirect handler can also tell you that you *shouldn't* update your address book.

### Example 11.13. Using the redirect handler to detect temporary redirects





- 
- ❶ This is a sample URL I've set up that is configured to tell clients to *temporarily* redirect to `http://www.example.com/redirect/`.
  - ❷ The server sends back a `302` status code, indicating a temporary redirect. The temporary new location of the data is given in the `Location` header.
  - ❸ `requests.get()` calls your `get` method, which calls the ancestor method of the same name in `requests.Session` which follows the redirect to the new location. Then your `get` method returns the response object.
  - ❹ And here you are, having successfully followed the redirect to `http://www.example.com/`. `response.status_code` tells you that this was a temporary redirect, which means that you should follow it.

## 11.8. Handling compressed data

The last important HTTP feature you want to support is compression. Many web services have the ability to send data compressed, which can save a lot of bandwidth.

Servers won't give you compressed data unless you tell them you can handle it.

### Example 11.14. Telling the server you would like compressed data





- 1 This is the key: once you've created your `obj` object, add an `header` to tell the server you can accept gzip-encoded data. `obj` is the name of the object.
- 2 There's your header going across the wire.
- 3 And here's what the server sends back: the `header` means that the data you're about to receive has been gzip-compressed.
- 4 The `header` is the length of the compressed data, not the uncompressed data. As you'll see in a minute, the actual length of the uncompressed data is much larger.

### Example 11.15. Decompressing the data



- 1 Continuing from the previous example, `f` is the file-like object returned from the URL opener. Using its `read()` method would ordinarily get you the data.
- 2 OK, this step is a little bit of messy workaround. Python has a `gzip` module, which reads (and actually writes) gzip-compressed files on disk.
- 3 Now you can create an instance of `GzipFile` and tell it that its "file" is the file-like object `f`.
- 4 This is the line that does all the actual work: "reading" from `g` will decompress the data. Strange? Yes, but it makes sense in a twisted kind of way.
- 5 Look ma, real data. (15955 bytes of it, in fact.)

"But wait!" I hear you cry. "This could be even easier!" I know what you're thinking. You're thinking that `r` returns a file-like object, so why not

### Example 11.16. Decompressing the data directly from the server

```
1 from io import BytesIO
2 import gzip
3
4 def get(url, headers=None, data=None, timeout=30):
5 r = requests.get(url, headers=headers, data=data, timeout=timeout)
6 if r.headers.get('Content-Encoding') == 'gzip':
7 buf = BytesIO(r.content)
8 f = gzip.GzipFile(fileobj=buf)
9 return f.read()
10 return r.content
```

- 1 Continuing from the previous example, you already have a `Response` object set up with an `Header`.
- 2 Simply opening the request will get you the headers (though not download any data yet). As you can see from the returned `Header`, this is a gzipped response.
- 3 Since `r` returns a file-like object, and you know from the headers that when you read it, you're going to get gzip-compressed data, why not

## 11.9. Putting it all together

You've seen all the pieces for building an intelligent HTTP web services client. Now let's see how they all fit together.

### Example 11.17. The `urllib2` function

This function is defined in `urllib2.py`

```
1 import urllib2
2 import sys
3 import re
4 import os
5 import socket
6 import ssl
7 import logging
8 import warnings
9 import errno
10 import time
11 import random
12 import hashlib
13 import mimetypes
14 import tempfile
15 import shutil
16 import subprocess
17 import ctypes
18 import ctypes.wintypes
19 import ctypes.util
20 import ctypes._endian
21 import ctypes._winapi
22 import ctypes._wintypes
23 import ctypes._wintypes
24 import ctypes._wintypes
25 import ctypes._wintypes
26 import ctypes._wintypes
27 import ctypes._wintypes
28 import ctypes._wintypes
29 import ctypes._wintypes
30 import ctypes._wintypes
31 import ctypes._wintypes
32 import ctypes._wintypes
33 import ctypes._wintypes
34 import ctypes._wintypes
35 import ctypes._wintypes
36 import ctypes._wintypes
37 import ctypes._wintypes
38 import ctypes._wintypes
39 import ctypes._wintypes
40 import ctypes._wintypes
41 import ctypes._wintypes
42 import ctypes._wintypes
43 import ctypes._wintypes
44 import ctypes._wintypes
45 import ctypes._wintypes
46 import ctypes._wintypes
47 import ctypes._wintypes
48 import ctypes._wintypes
49 import ctypes._wintypes
50 import ctypes._wintypes
51 import ctypes._wintypes
52 import ctypes._wintypes
53 import ctypes._wintypes
54 import ctypes._wintypes
55 import ctypes._wintypes
56 import ctypes._wintypes
57 import ctypes._wintypes
58 import ctypes._wintypes
59 import ctypes._wintypes
60 import ctypes._wintypes
61 import ctypes._wintypes
62 import ctypes._wintypes
63 import ctypes._wintypes
64 import ctypes._wintypes
65 import ctypes._wintypes
66 import ctypes._wintypes
67 import ctypes._wintypes
68 import ctypes._wintypes
69 import ctypes._wintypes
70 import ctypes._wintypes
71 import ctypes._wintypes
72 import ctypes._wintypes
73 import ctypes._wintypes
74 import ctypes._wintypes
75 import ctypes._wintypes
76 import ctypes._wintypes
77 import ctypes._wintypes
78 import ctypes._wintypes
79 import ctypes._wintypes
80 import ctypes._wintypes
81 import ctypes._wintypes
82 import ctypes._wintypes
83 import ctypes._wintypes
84 import ctypes._wintypes
85 import ctypes._wintypes
86 import ctypes._wintypes
87 import ctypes._wintypes
88 import ctypes._wintypes
89 import ctypes._wintypes
90 import ctypes._wintypes
91 import ctypes._wintypes
92 import ctypes._wintypes
93 import ctypes._wintypes
94 import ctypes._wintypes
95 import ctypes._wintypes
96 import ctypes._wintypes
97 import ctypes._wintypes
98 import ctypes._wintypes
99 import ctypes._wintypes
100 import ctypes._wintypes
```

- 1 `urllib2` is a handy utility module for, you guessed it, parsing URLs. Its primary function, also called `urllib2.urlopen`, takes a URL and splits it into a tuple of (scheme, netloc, path, query, fragment).
- 2 You identify yourself to the HTTP server with the `headers` passed in by the calling function. If no `headers` was specified, you use a default one defined in `urllib2`.
- 3 If an `auth` was given, send it in the `Header`.
- 4 If a last-modified date was given, send it in the `Header`.
- 5 Tell the server you would like compressed data if possible.
- 6 Build a URL opener that uses *both* of the custom URL handlers: `urllib2.HTTPHandler` for handling `GET` and `POST` requests, and `urllib2.HTTPSHandler` for handling `SSL` and other error conditions.

- 7 That's it! Open the URL and return a file-like object to the caller.

### Example 11.18. The `urlretrieve` function

This function is defined in `urllib`:

```
def urlretrieve(url, filename, reporthook=None, proxies=None):
```

- 1 First, you call the `urlretrieve` function with a URL, `hash`, `date`, and `proxies`.
- 2 Read the actual data returned from the server. This may be compressed; if so, you'll decompress it later.
- 3 Save the `hash` returned from the server, so the calling application can pass it back to you next time, and you can pass it on to `urlretrieve` which can
- 4 Save the `date` too.
- 5 If the server says that it sent compressed data, decompress it.
- 6 If you got a URL back from the server, save it, and assume that the status code is `200` until you find out otherwise.
- 7 If one of the custom URL handlers captured a status code, then save that too.

### Example 11.19. Using `urlretrieve`

```
urlretrieve(url, filename, reporthook=None, proxies=None)
```



- ❶ The very first time you fetch a resource, you don't have an `if_modified_since` or `if_none_match`, so you'll leave those out. (They're optional parameters.)
- ❷ What you get back is a dictionary of several useful headers, the HTTP status code, and the actual data returned from the server. `handle_redirect`
- ❸ If you ever get a `301` status code, that's a permanent redirect, and you need to update your URL to the new address.
- ❹ The second time you fetch the same resource, you have all sorts of information to pass back: a (possibly updated) URL, the `if_modified_since` from the first fetch, and the `if_none_match` from the first fetch.
- ❺ What you get back is again a dictionary, but the data hasn't changed, so all you got was a `304` status code and no data.

## 11.10. Summary

The `requests` module and its functions should now make perfect sense.

There are 5 important features of HTTP web services that every client should support:

- Identifying your application by setting a proper `User-Agent` header.
- Handling permanent redirects properly.
- Supporting `ETag` checking to avoid re-downloading data that hasn't changed.
- Supporting `Cache-Control` to avoid re-downloading data that hasn't changed.
- Supporting gzip compression to reduce bandwidth even when data *has* changed.

## Chapter 12. SOAP Web Services

Chapter 11 focused on document-oriented web services over HTTP. The "input parameter" was the URL, and the "return value" was an actual document.

This chapter will focus on SOAP web services, which take a more structured approach. Rather than dealing with HTTP requests and XML documents, SOAP uses a more structured approach.

SOAP is a complex specification, and it is somewhat misleading to say that SOAP is all about calling remote functions. Some people would prefer to think of it as a way to communicate with remote services.

### 12.1. Diving In

You use Google, right? It's a popular search engine. Have you ever wished you could programmatically access Google search results? Now you can.

#### Example 12.1.

```
#!/usr/bin/env python
```

```
import sys
```

```
import urllib
```

```
import urllib2
```

```
import re
```

```
def main():
```

```
 url = 'http://www.google.com/search?q=' + sys.argv[1] + '&btnG=Google'
```

```
 response = urllib2.urlopen(url)
```

```
 data = response.read()
```

```
 results = re.findall('<div class="search-result">', data)
```

```
 for result in results:
```

```
 print result
```

You can import this as a module and use it from a larger program, or you can run the script from the command line. On the command line, you can pass a search term as an argument.

Here is the sample output for a search for the word "python".

#### Example 12.2. Sample Usage of

```
python search.py python
```



### Further Reading on SOAP

<http://www.xmethods.net/> is a repository of public access SOAP web services.

The SOAP specification (<http://www.w3.org/TR/soap/>) is surprisingly readable, if you like that sort of thing.

## 12.2. Installing the SOAP Libraries

Unlike the other code in this book, this chapter relies on libraries that do not come pre-installed with Python.


Before you can dive into SOAP web services, you'll need to install three libraries: PyXML, fpconst, and SOAPpy.

### 12.2.1. Installing PyXML

The first library you need is PyXML, an advanced set of XML libraries that provide more functionality than the built-in XML libraries we st

#### Procedure 12.1.

Here is the procedure for installing PyXML:

1.  
Go to <http://pyxml.sourceforge.net/>, click Downloads, and download the latest version for your operating system.
2.  
If you are using Windows, there are several choices. Make sure to download the version of PyXML that matches the version of Pyth
3.  
Double-click the installer. If you download PyXML 0.8.3 for Windows and Python 2.3, the installer program will be 
4.  
Step through the installer program.
5.  
After the installation is complete, close the installer. There will not be any visible indication of success (no programs installed on the

To verify that you installed PyXML correctly, run your Python IDE and check the version of the XML libraries you have installed, as shown

### Example 12.3. Verifying PyXML Installation



This version number should match the version number of the PyXML installer program you downloaded and ran.

## 12.2.2. Installing fpconst

The second library you need is fpconst, a set of constants and functions for working with IEEE754 double-precision special values. This provides

### Procedure 12.2.

Here is the procedure for installing fpconst:

1.  
Download the latest version of fpconst from <http://www.analytics.washington.edu/statcomp/projects/rzope/fpconst/>.
2.  
There are two downloads available, one in `tar` format, the other in `zip` format. If you are using Windows, download the `zip` file; otherwise, do
3.  
Decompress the downloaded file. On Windows XP, you can right-click on the file and choose Extract All; on earlier versions of Windows
4.  
Open a command prompt and navigate to the directory where you decompressed the fpconst files.
5.  
Type `python setup.py install` to run the installation program.

To verify that you installed fpconst correctly, run your Python IDE and check the version number.

### Example 12.4. Verifying fpconst Installation



This version number should match the version number of the fpconst archive you downloaded and installed.

## 12.2.3. Installing SOAPpy

The third and final requirement is the SOAP library itself: SOAPpy.

### Procedure 12.3.

Here is the procedure for installing SOAPpy:



Go to <http://pywebsvcs.sourceforge.net/> and select Latest Official Release under the SOAPpy section.

2.  
There are two downloads available. If you are using Windows, download the `win32` file; otherwise, download the `unix` file.
3.  
Decompress the downloaded file, just as you did with `fpconst`.
4.  
Open a command prompt and navigate to the directory where you decompressed the SOAPpy files.
5.  
Type `python setup.py install` to run the installation program.

To verify that you installed SOAPpy correctly, run your Python IDE and check the version number.

### Example 12.5. Verifying SOAPpy Installation

```
Python 2.5.2 Shell
> import SOAPpy
> print SOAPpy.__version__
0.11.1
```

This version number should match the version number of the SOAPpy archive you downloaded and installed.

## 12.3. First Steps with SOAP

The heart of SOAP is the ability to call remote functions. There are a number of public access SOAP servers that provide simple functions for

The most popular public access SOAP server is <http://www.xmethods.net/>. This example uses a demonstration function that takes a United States

### Example 12.6. Getting the Current Temperature

```
Python 2.5.2 Shell
> from SOAPpy import *
> proxy = SOAPProxy('http://www.xmethods.net/demos/GetTemp')
> proxy.GetTemp('New York')
10
```

- ❶ You access the remote SOAP server through a proxy class, `SOAPProxy`. The proxy handles all the internals of SOAP for you, including creating the SOAP message.
- ❷ Every SOAP service has a URL which handles all the requests. The same URL is used for all function calls. This particular service on <http://www.xmethods.net/> is `http://www.xmethods.net/demos/GetTemp`.
- ❸ You're creating the `proxy` with the service URL and the service namespace. This doesn't make any connection to the SOAP server; it simply sets up the proxy.
- ❹ Now with everything configured properly, you can actually call remote SOAP methods as if they were local functions. You pass arguments to the method call.

Let's peek under those covers.

## 12.4. Debugging SOAP Web Services

The SOAP libraries provide an easy way to see what's going on behind the scenes.

Turning on debugging is a simple matter of setting two flags in the `ServiceProxy` configuration.

### Example 12.7. Debugging SOAP Web Services

```
1 from suds.client import ServiceProxy
2
3 # Create the proxy like normal, with the service URL and the namespace.
4 proxy = ServiceProxy('http://www.example.com/soap', namespace='http://www.example.com/soap')
5
6 # Turn on debugging by setting the following flags.
7 proxy.set_debug_options(SudsOptions.DEBUG_REQUEST, SudsOptions.DEBUG_RESPONSE)
8
9 # Call the remote SOAP method as usual. The SOAP library will print out both the outgoing XML request document, and the incoming XML response document.
```

- 1 First, create the `ServiceProxy` like normal, with the service URL and the namespace.
- 2 Second, turn on debugging by setting `SudsOptions.DEBUG_REQUEST` and `SudsOptions.DEBUG_RESPONSE`.
- 3 Third, call the remote SOAP method as usual. The SOAP library will print out both the outgoing XML request document, and the incoming XML response document.

Most of the XML request document that gets sent to the server is just boilerplate. Ignore all the namespace declarations; they're going to be t



- 1 The element name is the function name, `getWeather` is a dispatcher. Instead of calling separate local methods based on the method name, i
- 2 The function's XML element is contained in a specific namespace, which is the namespace you specified when you created the `Object`
- 3 The arguments of the function also got translated into XML. `introspects` each argument to determine its datatype (in this case it's a str

The XML return document is equally easy to understand, once you know what to ignore. Focus on this fragment within the `<Body>`



- 1 The server wraps the function return value within a `<Body>` element. By convention, this wrapper element is the name of the function, plus `<Body>`
- 2 The server returns the response in the same namespace we used in the request, the same namespace we specified when we first create
- 3 The return value is specified, along with its datatype (it's a float). `uses` this explicit datatype to create a Python object of the correct na

## 12.5. Introducing WSDL

The `Class` proxies local method calls and transparently turns them into invocations of remote SOAP methods. As you've seen, this is a lot of v

Consider this: the previous two sections showed an example of calling a simple remote SOAP method with one argument and one return value

That shouldn't come as a big surprise. If I wanted to call a local function, I would need to know what package or module it was in (the equiv

The big difference is introspection. As you saw in Chapter 4, Python excels at letting you discover things about modules and functions at run

WSDL lets you do that with SOAP web services. WSDL stands for "Web Services Description Language". Although designed to be flexible

A WSDL file is just that: a file. More specifically, it's an XML file. It usually lives on the same server you use to access the SOAP web servi

A WSDL file contains a description of everything involved in calling a SOAP web service:

- The service URL and namespace
- The type of web service (probably function calls using SOAP, although as I mentioned, WSDL is flexible enough to describe a wide
- The list of available functions
- The arguments for each function
- The datatype of each argument
- The return values of each function, and the datatype of each return value

In other words, a WSDL file tells you everything you need to know to be able to call a SOAP web service.

## 12.6. Introspecting SOAP Web Services with WSDL

Like many things in the web services arena, WSDL has a long and checkered history, full of political strife and intrigue. I will skip over this for now.

The most fundamental thing that WSDL allows you to do is discover the available methods offered by a SOAP server.

### Example 12.8. Discovering The Available Methods



- 1 SOAPpy includes a WSDL parser. At the time of this writing, it was labeled as being in the early stages of development, but I had no problem using it.
- 2 To use a WSDL file, you again use a proxy class, `SOAPProxy`, which takes a single argument: the WSDL file. Note that in this case you are passing the filename, not the URL.
- 3 The WSDL proxy class exposes the available functions as a Python dictionary, `methods`. So getting the list of available methods is as simple as `proxy.methods.keys()`.

Okay, so you know that this SOAP server offers a single method: `getStockPrice`. But how do you call it? The WSDL proxy object can tell you that too.

### Example 12.9. Discovering A Method's Arguments



- 1 The `methods` dictionary is filled with a SOAPpy-specific structure called `SOAPFunction`. A `SOAPFunction` object contains information about one specific function, including its name, arguments, and return type.
- 2 The function arguments are stored in `args`, which is a Python list of `SOAPArgument` objects that hold information about each parameter.
- 3 Each `SOAPArgument` object contains a `name` attribute, which is the argument name. You are not required to know the argument name to call the function though.
- 4 Each parameter is also explicitly typed, using datatypes defined in XML Schema. You saw this in the wire trace in the previous section.

WSDL also lets you introspect into a function's return values.

### Example 12.10. Discovering A Method's Return Values



- 1 The adjunct to `for` function arguments is `for` return value. It is also a list, because functions called through SOAP can return multiple values.
- 2 Each `Object` contains `and`. This function returns a single value, named `value` which is a float.

Let's put it all together, and call a SOAP web service through a WSDL proxy.

### Example 12.11. Calling A Web Service Through A WSDL Proxy


























































































































































- 1 The configuration is simpler than calling the SOAP service directly, since the WSDL file contains the both service URL and namespace.
- 2 Once the `Proxy` object is created, you can call a function as easily as you did with the `Service` object. This is not surprising; the `Proxy` is just a wrapper around the `Service`.
- 3 You can access the `Proxy` with `getProxy()`. This is useful to turning on debugging, so that when you can call functions through the WSDL proxy, its

## 12.7. Searching Google

Let's finally turn to the sample code that you saw at the beginning of this chapter, which does something more useful and exciting than get

Google provides a SOAP API for programmatically accessing Google search results. To use it, you will need to sign up for Google Web Services.

### Procedure 12.4. Signing Up for Google Web Services

1. Go to <http://www.google.com/apis/> and create a Google account. This requires only an email address. After you sign up you will receive a developer key.
2. Also on <http://www.google.com/apis/>, download the Google Web APIs developer kit. This includes some sample code in several programming languages.
3. Decompress the developer kit file and find `GoogleWebAPIs.jar`. Copy this file to some permanent location on your local drive. You will need it later in this chapter.

Once you have your developer key and your Google WSDL file in a known place, you can start poking around with Google Web Services.

### Example 12.12. Introspecting Google Web Services



- 1 Getting started with Google web services is easy: just create a `GoogleWebServices` object and point it at your local copy of Google's WSDL file.
- 2 According to the WSDL file, Google offers three functions: `search`, `searchCache`, and `searchIndex`. These do exactly what they sound like: perform a Google search and return the results from the cache, or return the results from the index.
- 3 The `search` function takes a number of parameters of various types. Note that while the WSDL file can tell you what the arguments are called, it doesn't tell you what they are used for.

Here is a brief synopsis of all the parameters to the `search` function:

- `apiKey`: Your Google API key, which you received when you signed up for Google web services.
- `q`: The search word or phrase you're looking for. The syntax is exactly the same as Google's web form, so if you know any advanced search techniques, you can use them here.
- `start`: The index of the result to start on. Like the interactive web version of Google, this function returns 10 results at a time. If you want to see the first 10 results, set this to 0.
- `numResults`: The number of results to return. Currently capped at 10, although you can specify fewer if you are only interested in a few results.
- `filter`: If `filter` is `True`, Google will filter out duplicate pages from the results.
- `country`: Set this to `plus a country code` to get results only from a particular country. Example: `GB` to search pages in the United Kingdom. You can also use `ALL` to search all countries.

- If `safe` Google will filter out porn sites.
- `lr` ("language restrict") - Set this to a language code to get results only in a particular language.
- `ie` and `oe` ("input encoding" and "output encoding") - Deprecated, both must be `utf-8`

### Example 12.13. Searching Google

```

1 from urllib2 import urlopen
2 import urllib
3 import sys
4
5 # Set up the search parameters
6 url = "http://www.google.com/search?q=%s" % sys.argv[1]
7
8 # Open the URL
9 f = urlopen(url)
10
11 # Read the response
12 data = f.read()
13
14 # Print the response
15 print data

```

- 1 After setting up the `urllib2` object, you can call `urlopen` with all ten parameters. Remember to use your own Google API key that you received when
- 2 There's a lot of information returned, but let's look at the actual search results first. They're stored in `data` and you can access them just like
- 3 Each element in the `data` is an object that has a `dict` and other useful attributes. At this point you can use normal Python introspection techn

The `urllib2` object contains more than the actual search results. It also contains information about the search itself, such as how long it took and how

### Example 12.14. Accessing Secondary Information From Google

```

1 from urllib2 import urlopen
2 import urllib
3 import sys
4
5 # Set up the search parameters
6 url = "http://www.google.com/search?q=%s" % sys.argv[1]
7
8 # Open the URL
9 f = urlopen(url)
10
11 # Read the response
12 data = f.read()
13
14 # Print the response
15 print data

```

- 1 This search took 0.224919 seconds. That does not include the time spent sending and receiving the actual SOAP XML documents. It's
- 2 In total, there were approximately 30 million results. You can access them 10 at a time by changing the `num` parameter and calling `urlopen` again.
- 3 For some queries, Google also returns a list of related categories in the Google Directory (<http://directory.google.com/>). You can appe

## 12.8. Troubleshooting SOAP Web Services

Of course, the world of SOAP web services is not all happiness and light. Sometimes things go wrong.

As you've seen throughout this chapter, SOAP involves several layers. There's the HTTP layer, since SOAP is sending XML documents to,

Beyond the underlying HTTP layer, there are a number of things that can go wrong. SOAPpy does an admirable job hiding the SOAP syntax

Here are a few examples of common mistakes that I've made in using SOAP web services, and the errors they generated.

### Example 12.15. Calling a Method With an Incorrectly Configured Proxy

```
#!/usr/bin/python
import sys
import SOAPpy
url = 'http://www.example.com/soap'
s = SOAPpy.SOAPProxy(url)
s.someMethod(1)
```

- ❶ Did you spot the mistake? You're creating a `SOAPProxy` manually, and you've correctly specified the service URL, but you haven't specified the proxy.
- ❷ The server responds by sending a SOAP Fault, which SOAPpy turns into a Python exception of type `SOAPpy.SAPFault`. All errors returned from any SOAP service are of this type.

Misconfiguring the basic elements of the SOAP service is one of the problems that WSDL aims to solve. The WSDL file contains the service

### Example 12.16. Calling a Method With the Wrong Arguments

```
#!/usr/bin/python
import sys
import SOAPpy
url = 'http://www.example.com/soap'
s = SOAPpy.SOAPProxy(url)
s.someMethod(1)
```

- ❶ Did you spot the mistake? It's a subtle one: you're calling `someMethod` with an integer instead of a string. As you saw from introspecting the WSDL file, the method expects a string.
- ❷ Again, the server returns a SOAP Fault, and the human-readable part of the error gives a clue as to the problem: you're calling a function with the wrong arguments.



It's also possible to write Python code that expects a different number of return values than the remote function actually returns.

### Example 12.17. Calling a Method and Expecting the Wrong Number of Return Values

南華真經

❶ Did you spot the mistake? `time` only returns one value, a float, but you've written code that assumes you're getting two values and trying to

What about Google's web service? The most common problem I've had with it is that I forget to set the application key properly.

### Example 12.18. Calling a Method With An Application-Specific Error





























- ❶ Can you spot the mistake? There's nothing wrong with the calling syntax, or the number of arguments, or the datatypes. The problem
- ❷ The Google server responds with a SOAP Fault and an incredibly long error message, which includes a complete Java stack trace. Ren

### Further Reading on Troubleshooting SOAP

New developments for SOAPpy (<http://www-106.ibm.com/developerworks/webservices/library/ws-pyth17.html>) steps through trying

## 12.9. Summary

SOAP web services are very complicated. The specification is very ambitious and tries to cover many different use cases for web services. Th

Before diving into the next chapter, make sure you're comfortable doing all of these things:

- Connecting to a SOAP server and calling remote methods.
- Loading a WSDL file and introspecting remote methods.
- Debugging SOAP calls with wire traces.
- Troubleshooting common SOAP-related errors.

# Chapter 13. Unit Testing

## 13.1. Introduction to Roman numerals

In previous chapters, you "dived in" by immediately looking at code and trying to understand it as quickly as possible. Now that you have some

In the next few chapters, you're going to write, debug, and optimize a set of utility functions to convert to and from Roman numerals. You saw

The rules for Roman numerals lead to a number of interesting observations:

There is only one correct way to represent a particular number as Roman numerals.1.

The converse is also true: if a string of characters is a valid Roman numeral, it represents only one number (*i.e.* it can only be read one

There is a limited range of numbers that can be expressed as Roman numerals, specifically 1 through 9. (The Romans did have several

There is no way to represent 0 in Roman numerals. (Amazingly, the ancient Romans had no concept of 0 as a number. Numbers were 1

There is no way to represent negative numbers in Roman numerals.5.

There is no way to represent fractions or non-integer numbers in Roman numerals.6.

Given all of this, what would you expect out of a set of functions to convert to and from Roman numerals?

### Requirements

It should return the Roman numeral representation for all integers 1 to 9.

It should fail when given an integer outside the range 1 to 9.

It should fail when given a non-integer number.3.

It should take a valid Roman numeral and return the number that it represents.4.

It should fail when given an invalid Roman numeral.5.

If you take a number, convert it to Roman numerals, then convert that back to a number, you should end up with the number you started

It should always return a Roman numeral using uppercase letters.7.

It should only accept uppercase Roman numerals (*i.e.* it should fail when given lowercase input).8.

### Further reading

This site (<http://www.wilkiecollins.demon.co.uk/roman/front.htm>) has more on Roman numerals, including a fascinating history (http://

## 13.2. Diving in

Now that you've completely defined the behavior you expect from your conversion functions, you're going to do something a little unexpected

This is called unit testing, since the set of two conversion functions can be written and tested as a unit, separate from any larger program they

is included with Python 2.1 and later. Python 2.0 users can download it from <http://pyunit.sourceforge.net/>.

Unit testing is an important part of an overall testing-centric development strategy. If you write unit tests, it is important to write them early (

Before writing code, it forces you to detail your requirements in a useful fashion.


While writing code, it keeps you from over-coding. When all the test cases pass, the function is complete.

When refactoring code, it assures you that the new version behaves the same way as the old version.

When maintaining code, it helps you cover your ass when someone comes screaming that your latest change broke their old code. ("

When writing code in a team, it increases confidence that the code you're about to commit isn't going to break other peoples' code,

## 13.3. Introducing

This is the complete test suite for your Roman numeral conversion functions, which are yet to be written but will eventually be in . It is in

### Example 13.1.

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5>).

```















































```





### Further reading

The PyUnit home page (<http://pyunit.sourceforge.net/>) has an in-depth discussion of using the framework (<http://pyunit.sourceforge.net/pyunit.html>) explains why test cases are stored separately (<http://pyunit.sourceforge.net/pyunit.html>) summarizes the <http://www.python.org/doc/current/lib/module-PyUnit.html> *Python Library Reference* (<http://www.python.org/doc/current/lib/>) discusses why you should write unit tests (<http://www.extremeprogramming.org/>) The Portland Pattern Repository (<http://www.c2.com/cgi/wiki/>) has an ongoing discussion of unit tests (<http://www.c2.com/cgi/wiki/>)

## 13.4. Testing for success

The most fundamental part of unit testing is constructing individual test cases. A test case answers a single question about the code it is testing.

A test case should be able to...

- ...run completely by itself, without any human input. Unit testing is about automation.
- ...determine by itself whether the function it is testing has passed or failed, without a human interpreting the results.
- ...run in isolation, separate from any other test cases (even if they test the same functions). Each test case is an island.

Given that, let's build the first test case. You have the following requirement:

It should return the Roman numeral representation for all integers from 1 to 9.

**Example 13.2.**

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99  
100





- ❶ To write a test case, first subclass the `TestCase` class of the `unittest` module. This class provides many useful methods which you can use in your test cases.
- ❷ This is a list of integer/numeral pairs that I verified manually. It includes the lowest ten numbers, the highest number, every number between 10 and 20, and so on.
- ❸ Every individual test is its own method, which must take no parameters and return no value. If the method exits normally without raising an exception, the test passes.
- ❹ Here you call the actual `isPrime` function. (Well, the function hasn't been written yet, but once it is, this is the line that will call it.) Notice that you are not trapping any exceptions when you call `isPrime`. This is intentional. `isPrime` shouldn't raise an exception when you call it with a valid integer.
- ❺ Also notice that you are not trapping any exceptions when you call `isPrime`. This is intentional. `isPrime` shouldn't raise an exception when you call it with a valid integer.
- ❻ Assuming the `isPrime` function was defined correctly, called correctly, completed successfully, and returned a value, the last step is to check whether the value is `True` or `False`.

## 13.5. Testing for failure

It is not enough to test that functions succeed when given good input; you must also test that they fail when given bad input. And not just any bad input.

Remember the other requirements for `isPrime`:

- ❶ `isPrime` should fail when given an integer outside the range 1 to 9.
- ❷ `isPrime` should fail when given a non-integer number.

In Python, functions indicate failure by raising exceptions, and the `unittest` module provides methods for testing whether a function raises a particular exception.

### Example 13.3. Testing bad input to `isPrime`



- 1 The `TestClass` of the `Test` provides the `test` method, which takes the following arguments: the exception you're expecting, the function you're testing
- 2 Along with testing numbers that are too large, you need to test numbers that are too small. Remember, Roman numerals cannot express
- 3 Requirement #3 specifies that `test` cannot accept a non-integer number, so here you test to make sure that `test` raises a `ValueError` when called

The next two requirements are similar to the first three, except they apply to  $\mathcal{A}$  instead of  $\mathcal{B}$ .

It should take a valid Roman numeral and return the number that it represents.

It should fail when given an invalid Roman numeral.5.

Requirement #4 is handled in the same way as requirement #1, iterating through a sampling of known values and testing each in turn. Requir

### Example 13.4. Testing bad input to `atoi`

- ❶ Not much new to say about these; the pattern is exactly the same as the one you used to test bad input to `atoi`. I will briefly note that you have

### 13.6. Testing for sanity

Often, you will find that a unit of code contains a set of reciprocal functions, usually in the form of conversion functions where one converts

Consider this requirement:

If you take a number, convert it to Roman numerals, then convert that back to a number, you should end up with the number you started with.

### Example 13.5. Testing $H_0$ against $H_1$

2

- 1 You've seen the `range` function before, but here it is called with two arguments, which returns a list of integers starting at the first argument and ending at the second argument.
- 2 I just wanted to mention in passing that `is` is not a keyword in Python; here it's just a variable name like any other.
- 3 The actual testing logic here is straightforward: take a number (`n`), convert it to a Roman numeral (`roman`), then convert it back to a number (`int(roman)`).

The last two requirements are different from the others because they seem both arbitrary and trivial:

- ❶ should always return a Roman numeral using uppercase letters.<sup>7</sup>
- ❷ should only accept uppercase Roman numerals (*i.e.* it should fail when given lowercase input).<sup>8</sup>

In fact, they are somewhat arbitrary. You could, for instance, have stipulated that ❸ accept lowercase and mixed case input. But they are not co

### Example 13.6. Testing for case

```
❶ def to_roman(n):
❷ """Convert the number n to a Roman numeral.
❸ """
❹ if not isinstance(n, int):
❺ raise TypeError("non-int input")
❻ if n < 1 or n > 4000:
❼ raise ValueError("argument out of range")
❽ int_to_roman = {1: "I", 4: "IV", 5: "V", 9: "IX", 10: "X",
❾ 40: "XL", 50: "L", 90: "XC", 100: "C", 400: "CD", 500: "D", 900: "CM"}
❿ result = ""
⓫ while n > 0:
⓬ for (p, s) in sorted(int_to_roman.items(), reverse=True):
⓭ if n >= p:
⓮ result += s
⓯ n -= p
⓰ return result
```

- ❶ The main point of testing things about the function is not that the function should work, but only testing the question: "Imagine if you had combined the two tests into one, and the function returned uppercase. Would it still work?"
- ❷ There's a similar lesson to be learned here: even though "you know" that ❶ always returns uppercase, you are explicitly converting its return value to uppercase.
- ❸ Note that you're not assigning the return value of ❶ to anything. This is legal syntax in Python; if a function returns a value but nobody's interested in it, you can just ignore it.
- ❹ This is a complicated line, but it's very similar to what you did in the ❶ and ❷ tests. You are testing to make sure that calling a particular function with a particular argument returns a particular value.

In the next chapter, you'll see how to write code that passes these tests.

---

[6] I can resist everything except temptation." --Oscar Wilde

### 14.1. stage 1

### Example 14.1.

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5>).

- 1 This is how you define your own custom exceptions in Python. Exceptions are classes, and you create your own by subclassing existing ones.
- 2 The `InvalidInputException` will eventually be used by `validate_email` to flag various forms of invalid input, as specified in `validate_email.py`.
- 3 The `InvalidInputException` will eventually be used by `validate_email` to flag invalid input, as specified in `validate_email.py`.
- 4 At this stage, you want to define the API of each of your functions, but you don't want to code them yet, so you stub them out using the `__main__` block.

Run `test` with the `-v` command-line option, which will give more verbose output so you can see exactly what's going on as each test case runs. With

一、**《说文解字》**：中国第一部系统分析汉字字形、考究字源的字典，由东汉许慎编著。





- 1 Running the script runs `test.py` which runs each test case, which is to say each method defined in each class within `test.py`. For each test case, it prints the result.
- 2 For each failed test case, `test.py` displays the trace information showing exactly what happened. In this case, the call to `test()` also called `test()` raised a `ValueError`.
- 3 After the detail, `test.py` displays a summary of how many tests were performed and how long it took.
- 4 Overall, the unit test failed because at least one test case did not pass. When a test case doesn't pass, `test.py` distinguishes between failures and errors.

## 14.2. Stage 2

Now that you have the framework of the `test` module laid out, it's time to start writing code and passing test cases.

### Example 14.3. `test.py`

This file is available in `test.py` in the examples directory.

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5>).



```
1
2
3
```

- 1 `roman_numerals` is a tuple of tuples which defines three things:
  - The character representations of the most basic Roman numerals. Note that this is not just the single-character Roman numerals, but also the subtractive combinations like `IV` and `IX`.
  - The order of the Roman numerals. They are listed in descending value order, from `M` all the way down to `I`.
  - The value of each Roman numeral. Each inner tuple is a pair of `(character, value)`.
- 2 Here's where your rich data structure pays off, because you don't need any special logic to handle the subtraction rule. To convert to Roman, you just iterate over the numerals in order, and if the value is greater than or equal to the current numeral's value, you use that numeral and subtract its value from the input.

### Example 14.4. How `to_roman` works

If you're not clear how `to_roman` works, add a `print` statement to the end of the `while` loop:

```
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

So `to_roman` appears to work, at least in this manual spot check. But will it pass the unit testing? Well no, not entirely.

### Example 14.5. Output of `test_roman.py` against `to_roman`

Remember to run `test_roman.py` with the `-v` command-line flag to enable verbose mode.

```
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

- 1 `to_roman` does, in fact, always return uppercase, because `roman_numerals` defines the Roman numeral representations as uppercase. So this test passes already.
- 2 Here's the big news: this version of the `to_roman` function passes the known values test. Remember, it's not comprehensive, but it does put the test through its paces.
- 3 However, the function does not "work" for bad values; it fails every single bad input test. That makes sense, because you didn't include any logic to handle bad input.



Here's the rest of the output of the unit test, listing the details of all the failures. You're down to 10.

**一、新學級獎勵制度**

**二、新學級獎勵制度**

**三、新學級獎勵制度**

**四、新學級獎勵制度**

**五、新學級獎勵制度**

**六、新學級獎勵制度**

**七、新學級獎勵制度**

**八、新學級獎勵制度**

**九、新學級獎勵制度**

**十、新學級獎勵制度**

```
def is_int(s):
 """Return True if s is an integer, False otherwise.
 This function uses a try/except block to attempt to
 convert s to an integer. If successful, it returns
 True. If a ValueError is raised, it returns False.
 """
 try:
 int(s)
 except ValueError:
 return False
 return True

def is_float(s):
 """Return True if s is a float, False otherwise.
 This function uses a try/except block to attempt to
 convert s to a float. If successful, it returns
 True. If a ValueError is raised, it returns False.
 """
 try:
 float(s)
 except ValueError:
 return False
 return True

def is_int_or_float(s):
 """Return True if s is an integer or a float,
 False otherwise.
 """
 return is_int(s) or is_float(s)

def is_int_or_float_or_str(s):
 """Return True if s is an integer, a float, or a
 string, False otherwise.
 """
 return is_int(s) or is_float(s) or isinstance(s, str)

def is_int_or_float_or_str_or_bool(s):
 """Return True if s is an integer, a float, a
 string, or a boolean, False otherwise.
 """
 return is_int(s) or is_float(s) or isinstance(s, str) or \
 isinstance(s, bool)
```

### 14.3. stage 3

Now that `is_int` behaves correctly with good input (integers from `1` to `9`) it's time to make it behave correctly with bad input (everything else).

#### Example 14.6. `is_int`

This file is available in `04` in the examples directory.

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5>).

```
def is_int(s):
```

```
 try:
```

```
 int(s)
```

```
 except ValueError:
```

```
 return False
```

```
 else:
```

```
 return True
```

```
def to_roman(n):
```

```
 if not isinstance(n, int):
```

```
 raise ValueError('non-integer value')
```

```
 if n < 1 or n > 4999:
```

```
 raise ValueError('number out of range')
```

```
 if not is_int(n):
```

```
 raise ValueError('non-integer value')
```

```
 if n < 1:
```

```
 raise ValueError('number out of range')
```

```
 if n > 4999:
```

```
 raise ValueError('number out of range')
```

```
 if not is_int(n):
```

```
 raise ValueError('non-integer value')
```

```
 if n < 1:
```

```
 raise ValueError('number out of range')
```

```
 if n > 4999:
```

```
 raise ValueError('number out of range')
```

```
 if not is_int(n):
```

```
 raise ValueError('non-integer value')
```

```
 if n < 1:
```

```
 raise ValueError('number out of range')
```

```
 if n > 4999:
```

```
 raise ValueError('number out of range')
```

```
 if not is_int(n):
```

```
 raise ValueError('non-integer value')
```

```
 if n < 1:
```

```
 raise ValueError('number out of range')
```

```
 if n > 4999:
```

```
 raise ValueError('number out of range')
```

```
 if not is_int(n):
```

```
 raise ValueError('non-integer value')
```

```
 if n < 1:
```

```
 raise ValueError('number out of range')
```

```
 if n > 4999:
```

```
 raise ValueError('number out of range')
```

```
 if not is_int(n):
```

```
 raise ValueError('non-integer value')
```

```
 if n < 1:
```

```
 raise ValueError('number out of range')
```

```
 if n > 4999:
```

```
 raise ValueError('number out of range')
```

```
 if not is_int(n):
```

```
 raise ValueError('non-integer value')
```

- ❶ This is a nice Pythonic shortcut: multiple comparisons at once. This is equivalent to `n < 1 or n > 4999` but it's much easier to read. This is the range check.
- ❷ You raise exceptions yourself with the `raise` statement. You can raise any of the built-in exceptions, or you can raise any of your custom exceptions.
- ❸ This is the non-integer check. Non-integers can not be converted to Roman numerals.
- ❹ The rest of the function is unchanged.

## Example 14.7. Watching `try` handle bad input

```
def to_roman(n):
```

```
 try:
```

```
 int(s)
```

```
 except ValueError:
```

```
 return False
```

```
 else:
```

```
 return True
```



### Example 14.8. Output of `test_is_int` against `is_int`



- 1 `is_int` still passes the known values test, which is comforting. All the tests that passed in stage 2 still pass, so the latest code hasn't broken anything.
- 2 More exciting is the fact that all of the bad input tests now pass. This test, `test_is_int_bad_input`, passes because of the `isinstance` check. When a non-integer is passed to `is_int`, it returns `False`.
- 3 This test, `test_is_int_exception`, passes because of the `isinstance` check, which raises an `AssertionError`, which is what `test_is_int_exception` is looking for.





❶ You're down to 6 failures, and all of them involve the known values test, the three separate bad input tests, the case check, and the sa

The most important thing that comprehensive unit testing can tell you is when to stop coding. When all the unit tests for a function pass, stop

## 14.4. stage 4

Now that's done, it's time to start coding. Thanks to the rich data structure that maps individual Roman numerals to integer values, this is n

### Example 14.9.

This file is available in the examples directory.

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5>).

```
1
```

```
2
3
4
5
6
7
```

```
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

```
101
```

```
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
```

❶ The pattern here is the same as ❶. You iterate through your Roman numeral data structure (a tuple of tuples), and instead of matching the

### Example 14.10. How ❶ works

If you're not clear how ❶ works, add a `print` statement to the end of the `while` loop:

```
1
2
3
4
5
6
```

```
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```



### Example 14.11. Output of `isinstance` against `is`



- ❶ Two pieces of exciting news here. The first is that `isinstance` works for good input, at least for all the known values you test.
- ❷ The second is that the sanity check also passed. Combined with the known values tests, you can be reasonably sure that both `isinstance` and `is` work.



## 14.5. stage 5

Now that it works properly with good input, it's time to fit in the last piece of the puzzle: making it work properly with bad input. That means

If you're not familiar with regular expressions and didn't read Chapter 7, *Regular Expressions*, now would be a good time.

As you saw in Section 7.3, “Case Study: Roman Numerals”, there are several simple rules for constructing a Roman numeral, using the letter

Characters are additive.  $\mathbf{I}_s \mathbf{1} \mathbf{I}_s \mathbf{2}$  and  $\mathbf{I}_s \mathbf{3} \mathbf{I}_s \mathbf{6}$  literally, "5 and 1, 7 and 8.

The tens characters (IXC and M) can be repeated up to three times. At 4 you need to subtract from the next highest fives character. You

Similarly, at 9 you need to subtract from the next highest tens character: 8s  $\nabla$  but 9s  $\times$  less than 0, not  $\nabla$  since the  $\nabla$  character can not

The five characters can not be repeated. As always represented as X never as V As always C never L.

Roman numerals are always written highest to lowest, and read left to right, so order of characters matters very much. **Is 99** a com

### Example 14.12.

This file is available in `lib` in the examples directory.

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5>).





- 1 This is just a continuation of the pattern you discussed in Section 7.3, “Case Study: Roman Numerals”. The tens places is either `IX` or `X`.
- 2 Having encoded all that logic into a regular expression, the code to check for invalid Roman numerals becomes trivial. If `is_valid` returns an `o`

At this point, you are allowed to be skeptical that that big ugly regular expression could possibly catch all the types of invalid Roman numerals.

**Example 14.13. Output of `is_valid` against `roman`**



- ❶ One thing I didn't mention about regular expressions is that, by default, they are case-sensitive. Since the regular expression `was` expression
- ❷ More importantly, the bad input tests pass. For instance, the malformed antecedents test checks cases like `As you've seen`, this does not
- ❸ In fact, all the bad input tests pass. This regular expression catches everything you could think of when you made your test cases.
- ❹ And the anticlimax award of the year goes to the word `"` which is printed by the `module` when all the tests pass.

When all of your tests pass, stop coding.

# Chapter 15. Refactoring

## 15.1. Handling bugs

Despite your best efforts to write comprehensive unit tests, bugs happen. What do I mean by "bug"? A bug is a test case you haven't written y

### Example 15.1. The bug



- 1 Remember in the previous section when you kept seeing that an empty string would match the regular expression you were using to cl

After reproducing the bug, and before fixing it, you should write a test case that fails, thus illustrating the bug.

### Example 15.2. Testing for the bug (



- 1 Pretty simple stuff here. Call with an empty string and make sure it raises an exception. The hard part was finding the bug; now that y

Since your code has a bug, and you now have a test case that tests this bug, the test case will fail:

### Example 15.3. Output of against





Now you can fix the bug.

### Example 15.4. Fixing the bug (1)

This file is available in `01` in the examples directory.



- 1 Only two lines of code are required: an explicit check for an empty string, and a `return` statement.

### Example 15.5. Output of `test_is_palindrome` against `02`



- 1 The blank string test case now passes, so the bug is fixed.
- 2 All the other test cases still pass, which means that this bug fix didn't break anything else. Stop coding.

Coding this way does not make fixing bugs any easier. Simple bugs (like this one) require simple test cases; complex bugs will require complex test cases.

## 15.2. Handling changing requirements

Despite your best efforts to pin your customers to the ground and extract exact requirements from them on pain of horrible nasty things invol

Suppose, for instance, that you wanted to expand the range of the Roman numeral conversion functions. Remember the rule that said that no

### Example 15.6. Modifying test cases for new requirements (1)

This file is available in `lib` in the examples directory.

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5>).

[illegible]



[illegible]

- 1 The existing known values don't change (they're all still reasonable values to test), but you need to add a few more in the `range`. Here
- 2 The definition of "large input" has changed. This test used to call `h` with `0` and expect an error; now that `0` are good values, you need to bu

- ③ The definition of "too many repeated numerals" has also changed. This test used to call `is_valid` with `M` and expect an error; now that `M`s are considered valid, this test fails.
- ④ The sanity check and case checks loop through every number in the range, from `1` to `9`. Since the range has now expanded, these loops now fail.

Now your test cases are up to date with the new requirements, but your code is not, so you expect several of the test cases to fail.

### Example 15.7. Output of `pytest` against `test_roman.py`



- ① Our case checks now fail because they loop from `1` to `9` but `is_valid` only accepts numbers from `1` to `9`, so it will fail as soon as the test case hits `0`.
- ② The `known values` test will fail as soon as it hits `M` because `is_valid` still thinks this is an invalid Roman numeral.
- ③ The `known values` test will fail as soon as it hits `9` because `is_valid` still thinks this is out of range.
- ④ The sanity check will also fail as soon as it hits `9` because `is_valid` still thinks this is out of range.












● ● ● ● ● ●

卷之四











- 1 Only needs one small change, in the range check. Where you used to check `5` you now check `8`. And you change the error message that says "must be between 5 and 10" to "must be between 8 and 10".
- 2 You don't need to make any changes to `isInt` at all. The only change is to `isInRange`. If you look closely, you'll notice that you added another optional parameter to `isInRange`.

You may be skeptical that these two small changes are all that you need. Hey, don't take my word for it; see for yourself:

### Example 15.9. Output of `ls -la` against `/`












- 1 All the test cases pass. Stop coding.

Comprehensive unit testing means never having to rely on a programmer who says "Trust me."

## 15.3. Refactoring

The best thing about comprehensive unit testing is not the feeling you get when all your test cases finally pass, or even the feeling you get when

Refactoring is the process of taking working code and making it work better. Usually, "better" means "faster", although it can also mean "using

Here, "better" means "faster". Specifically, the `is_match` function is slower than it needs to be, because of that big nasty regular expression that you use

### Example 15.10. Compiling regular expressions



- 1 This is the syntax you've seen before: `re.match(pattern, string)` takes a regular expression as a string (`pattern`) and a string to match against it (`string`). If the pattern matches, it returns a match object; otherwise, it returns `None`.
- 2 This is the new syntax: `re.compile(pattern)` takes a regular expression as a string and returns a pattern object. Note there is no string to match here. Compiling a regular expression is a one-time cost.
- 3 The compiled pattern object returned from `re.compile` has several useful-looking functions, including several (like `match` and `search`) that are available directly on the pattern object.
- 4 Calling the compiled pattern object's `match` function with the string `string` accomplishes the same thing as calling `re.match` with both the regular expression and the string.

Whenever you are going to use a regular expression more than once, you should compile it to get a pattern object, then call the methods on the pattern object.

### Example 15.11. Compiled regular expressions in `re`

This file is available in `examples` in the examples directory.

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5.0.tar.gz>).

```
1 import re
2
3 # Compile the regular expression
4 re.compile(r'([a-zA-Z0-9_]+@[a-zA-Z0-9_]+\.[a-zA-Z0-9_]+\.[a-zA-Z0-9_]+)')
5
6 # Test the regular expression
7 re.compile(r'([a-zA-Z0-9_]+@[a-zA-Z0-9_]+\.[a-zA-Z0-9_]+\.[a-zA-Z0-9_]+)').search('foo@bar.baz.com')
```

- 1 This looks very similar, but in fact a lot has changed. `re.compile` is no longer a string; it is a pattern object which was returned from `re.compile`.
- 2 That means that you can call methods on `re.compile` directly. This will be much, much faster than calling `re.compile` every time. The regular expression is compiled once.

So how much faster is it to compile regular expressions? See for yourself:

### Example 15.12. Output of `python -m unittest` against `re`

```
1 -
2
3
```

- 1 Just a note in passing here: this time, I ran the unit test *without* the `-v` option, so instead of the full `Test` for each test, you only get a dot for each test.
- 2 You ran 3 tests in 0.000 seconds, compared to 0.000 seconds without precompiling the regular expressions. That's an improvement overall, and re.compile is faster.
- 3 Oh, and in case you were wondering, precompiling the regular expression didn't break anything, and you just proved it.

There is one other performance optimization that I want to try. Given the complexity of regular expression syntax, it should come as no surprise that

### Example 15.13. `re.compile`

This file is available in `examples` directory.

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5.0.tar.gz>).

```
1 import re
2
3 # Compile the regular expression
4 re.compile(r'([a-zA-Z0-9_]+@[a-zA-Z0-9_]+\.[a-zA-Z0-9_]+\.[a-zA-Z0-9_]+)')
5
6 # Test the regular expression
7 re.compile(r'([a-zA-Z0-9_]+@[a-zA-Z0-9_]+\.[a-zA-Z0-9_]+\.[a-zA-Z0-9_]+)').search('foo@bar.baz.com')
```





- ❶ This new, "verbose" version runs at exactly the same speed as the old version. In fact, the compiled pattern objects are the same, since
- ❷ This new, "verbose" version passes all the same tests as the old version. Nothing has changed, except that the programmer who comes

## 15.4. Postscript

A clever reader read the previous section and took it to the next level. The biggest headache (and performance drain) in the program as it is c  
And best of all, he already had a complete set of unit tests. He changed over half the code in the module, but the unit tests stayed the same, so

### Example 15.17. 🐍

This file is available in 📄 in the examples directory.

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5>).



```

1 #!/usr/bin/env python
2 """
3 A simple program to demonstrate the use of the sys module.
4 It prints out the command-line arguments.
5 """
6 import sys
7
8 # Print out the command-line arguments
9 for i in range(1, len(sys.argv)):
10 print(sys.argv[i])

```

So how fast is it?

### Example 15.18. Output of `python3.6.1.py` against `python3.6.1`

```

1 .
2 -
3 [
4 [

```

Remember, the best performance you ever got in the original version was 13 tests in 3.315 seconds. Of course, it's not entirely a fair comparison.

The moral of the story?

Simplicity is a virtue.

Especially when regular expressions are involved.

And unit tests can give you the confidence to do large-scale refactoring... even if you didn't write the original code.

## 15.5. Summary

Unit testing is a powerful concept which, if properly implemented, can both reduce maintenance costs and increase flexibility in any long-term

This chapter covered a lot of ground, and much of it wasn't even Python-specific. There are unit testing frameworks for many languages, all

- Designing test cases that are specific, automated, and independent.

- Writing test cases *before* the code they are testing.

- Writing tests that test good input and check for proper results.

- Writing tests that test bad input and check for proper failures.

- Writing and updating test cases to illustrate bugs or reflect new requirements.

- Refactoring mercilessly to improve performance, scalability, readability, maintainability, or whatever other -ility you're lacking.

Additionally, you should be comfortable doing all of the following Python-specific things:

- Subclassing `unittest.TestCase` and writing methods for individual test cases.

- Using `assertEqual` to check that a function returns a known value.

- Using `assertRaises` to check that a function raises a known exception.

- Calling `unittest.main()` in your `__main__` clause to run all your test cases at once.

- Running unit tests in verbose or regular mode.

### Further reading

XProgramming.com (<http://www.xprogramming.com/>) has links to download unit testing frameworks (<http://www.xprogramming.com/>)



## Chapter 16. Functional Programming

### 16.1. Diving in

In Chapter 13, *Unit Testing*, you learned about the philosophy of unit testing. In Chapter 14, *Test-First Programming*, you stepped through the

The following is a complete Python program that acts as a cheap and simple regression testing framework. It takes unit tests that you've written

#### Example 16.1.

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5.0.tar.gz>)

```
#!/usr/bin/env python
import sys
import os
import glob
import unittest

def main():
 """
 Run all unit tests in the current directory.
 """
 # Find all test files
 test_files = glob.glob('*.py')

 # Run the tests
 test_loader = unittest.TestLoader()
 test_suite = test_loader.discover('.')
 runner = unittest.TextTestRunner(verbosity=2)
 runner.run(test_suite)

if __name__ == '__main__':
 main()
```

Running this script in the same directory as the rest of the example scripts that come with this book will find all the unit tests, named `test_*.py`, and run them.

#### Example 16.2. Sample output of

```
test_1.py
test_2.py
test_3.py
test_4.py
test_5.py
test_6.py
test_7.py
test_8.py
test_9.py
test_10.py
```



- ❶ The first 5 tests are from `test_introspection.py` which tests the example script from Chapter 4, *The Power Of Introspection*.
- ❷ The next 5 tests are from `test_first_program.py` which tests the example script from Chapter 2, *Your First Python Program*.
- ❸ The rest are from `test_unit_testing.py` which you studied in depth in Chapter 13, *Unit Testing*.

## 16.2. Finding the path

When running Python scripts from the command line, it is sometimes useful to know where the currently running script is located on disk.

This is one of those obscure little tricks that is virtually impossible to figure out on your own, but simple to remember once you see it. The key

### Example 16.3. `__file__`

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5.0.tar.gz>).



- ❶ Regardless of how you run a script, `__file__` will always contain the name of the script, exactly as it appears on the command line. This may or may not be the full path.
- ❷ `os.path.dirname()` takes a filename as a string and returns the directory path portion. If the given filename does not include any path information, `os.path.dirname()` returns `.`.
- ❸ `os.path.abspath()` is the key here. It takes a pathname, which can be partial or even blank, and returns a fully qualified pathname.

deserves further explanation. It is very flexible; it can take any kind of pathname.

#### Example 16.4. Further explanation of `os.getcwd()`

```
1 #!/usr/bin/env python
2 """
3 Example 16.4: Further explanation of os.getcwd()
4 """
5 import os
6
7 # Get the current working directory
8 cwd = os.getcwd()
9
10 # Print it out
11 print(cwd)
```

- 1 `os.getcwd()` returns the current working directory.
- 2 Calling `os.getcwd()` with an empty string returns the current working directory, same as `os.getcwd()`.
- 3 Calling `os.getcwd()` with a partial pathname constructs a fully qualified pathname out of it, based on the current working directory.
- 4 Calling `os.getcwd()` with a full pathname simply returns it.
- 5 `os.getcwd()` also *normalizes* the pathname it returns. Note that this example worked even though I don't actually have a 'foo' directory. `os.getcwd()` never checks if the directory exists.

The pathnames and filenames you pass to `os.getcwd()` do not need to exist.

`os.getcwd()` not only constructs full path names, it also normalizes them. That means that if you are in the `foo` directory, `os.getcwd()` will return `/usr/bin/foo`. It normalizes the path.

#### Example 16.5. Sample output from `os.getcwd()`

```
1 #!/usr/bin/env python
2 """
3 Example 16.5: Sample output from os.getcwd()
4 """
5 import os
6
7 # Get the current working directory
8 cwd = os.getcwd()
9
10 # Print it out
11 print(cwd)
```

- 1 In the first case, `os.getcwd()` includes the full path of the script. You can then use the `os.path.dirname()` function to strip off the script name and return the full directory.
- 2 If the script is run by using a partial pathname, `os.getcwd()` will still contain exactly what appears on the command line. `os.getcwd()` will then give you a partial pathname.
- 3 If the script is run from the current directory without giving any path, `os.getcwd()` will simply return an empty string. Given an empty string, `os.getcwd()` will return the current working directory.

Like the other functions in the `os` module, `is_cross_platform` is cross-platform. Your results will look slightly different than my examples if you're running **Addendum**. One reader was dissatisfied with this solution, and wanted to be able to run all the unit tests in the current directory, not the directory where the script is located.

**Example 16.6. Running scripts in the current directory**

- ```
#!/usr/bin/env python
import os
import sys
import unittest

# Add the current directory to the Python library search path
sys.path.append(os.getcwd())

# Import the unit test modules
import test_foo
import test_bar

if __name__ == '__main__':
    unittest.main()
```
- 1 Instead of setting `sys.path` to the directory where the currently running script is located, you set it to the current working directory instead. This is done by calling `os.getcwd()`.
 - 2 Append this directory to the Python library search path, so that when you dynamically import the unit test modules later, Python can find them.
 - 3 The rest of the function is the same.

This technique will allow you to re-use this script on multiple projects. Just put the script in a common directory, then change to the project's directory and run the script.

16.3. Filtering lists revisited

You're already familiar with using list comprehensions to filter lists. There is another way to accomplish this same thing, which some people find more intuitive. The function `filter()` takes a function and a list as arguments and returns a list containing all the elements from the list for which the function returns a true value. The function `filter()` is defined in the `itertools` module.

Got all that? It's not as difficult as it sounds.

Example 16.7. Introducing `filter()`

- ```
#!/usr/bin/env python
import sys
import itertools

def is_even(n):
 return n % 2 == 0

def filter_even(list):
 return list(filter(is_even, list))

if __name__ == '__main__':
 list = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
 filtered_list = filter_even(list)
 print filtered_list
```
- 1 `is_even` uses the built-in `mod` function `%` to return `True` if `n` is odd and `False` if `n` is even.
  - 2 `filter_even` takes two arguments, a function (`is_even`) and a list (`list`). It loops through the list and calls `is_even` with each element. If `is_even` returns a true value (remember, `True` is the only value that is considered "true" in Python), `filter` includes it in the new list.
  - 3 You could accomplish the same thing using list comprehensions, as you saw in Section 4.5, "Filtering Lists".
  - 4 You could also accomplish the same thing with a `for` loop. Depending on your programming background, this may seem more "straightforward".

**Example 16.8. Using `filter()`**

```
#!/usr/bin/env python
import sys
import itertools
```

❶

- ❶ As you saw in Section 16.2, “Finding the path”, `__file__` may contain the full or partial pathname of the directory of the currently running script.
- ❷ This is a compiled regular expression. As you saw in Section 15.3, “Refactoring”, if you’re going to use the same regular expression often, it’s a good idea to compile it once and reuse the compiled object.
- ❸ For each element in the `files` list, you’re going to call the `match` method of the compiled regular expression object, `re_obj`. If the regular expression matches, `re_obj.match()` returns a `Match` object; otherwise, it returns `None`.

**Historical note.** Versions of Python prior to 2.0 did not have list comprehensions, so you couldn’t filter using list comprehensions; the `filter` function was used instead.

### Example 16.9. Filtering using list comprehensions instead

```
1
2
3
```

- ❶ This will accomplish exactly the same result as using the `filter` function. Which way is more expressive? That’s up to you.

## 16.4. Mapping lists revisited

You’re already familiar with using list comprehensions to map one list into another. There is another way to accomplish the same thing, using the `map` function.

### Example 16.10. Introducing `map`

```
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

- ❶ `map` is a function that takes a function and a list as arguments. It calls the function with each element of the list in order. In this case, the function simply multiplies each element of the list by 2.
- ❷ You could accomplish the same thing with a list comprehension. List comprehensions were first introduced in Python 2.0; `map` has been around since Python 1.0.
- ❸ You could, if you insist on thinking like a Visual Basic programmer, use a `for` loop to accomplish the same thing.

### Example 16.11. `map` with lists of mixed datatypes

```
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

- ❶ As a side note, I’d like to point out that `map` works just as well with lists of mixed datatypes, as long as the function you’re using correctly handles the different types.

All right, enough play time. Let’s look at some real code.

DD

- As

10

## By

In

Hc

No

Ins

Oh

I re

Le

10

## OK

Fin

9

- Nc
























- 1 The built-in `__import__` function accomplishes the same goal as using the `import` statement, but it's an actual function, and it takes a string as an argument.
- 2 The variable `sys` is now the `sys` module, just as if you had said `import sys`. The variable `os` is now the `os` module, and so forth.

So `__import__` imports a module, but takes a string argument to do it. In this case the module you imported was just a hard-coded string, but it could just as easily be a variable.

### Example 16.15. Importing a list of modules dynamically

```
1 # Import a list of modules dynamically
2 import sys
3 from os.path import join, dirname
4 from os import listdir
5 from glob import glob
6 from sys import argv
7 from sys import exit
8 from sys import stderr
9 from sys import stdout
10 from sys import stdin
```

- 1 `__import__` is just a list of strings. Nothing fancy, except that the strings happen to be names of modules that you could import, if you wanted to.
- 2 Surprise, you wanted to import them, and you did, by mapping the `__import__` function onto the list. Remember, this takes each element of the list and imports it.
- 3 So now from a list of strings, you've created a list of actual modules. (Your paths may be different, depending on your operating system.)
- 4 To drive home the point that these are real modules, let's look at some module attributes. Remember, `sys` is the `sys` module, so `sys.__name__` is `'sys'`. All the other modules have their own names.

Now you should be able to put this all together and figure out what most of this chapter's code sample is doing.

## 16.7. Putting it all together

You've learned enough now to deconstruct the first seven lines of this chapter's code sample: reading a directory and importing selected modules.

### Example 16.16. The `__import__` function

```
1 # Import a list of modules dynamically
2 import sys
3 from os.path import join, dirname
4 from os import listdir
5 from glob import glob
6 from sys import argv
7 from sys import exit
8 from sys import stderr
9 from sys import stdout
10 from sys import stdin
```

Let's look at it line by line, interactively. Assume that the current directory is `/`, which contains the examples that come with this book, including this one.

### Example 16.17. Step 1: Get all the files



- ❶ `ls` is a list of all the files and directories in the script's directory. (If you've been running some of the examples already, you may also see

### Example 16.18. Step 2: Filter to find the files you care about



- ❶ This regular expression will match any string that ends with `.py`. Note that you need to escape the period, since a period in a regular expression means "any character".
- ❷ The compiled regular expression acts like a function, so you can use it to filter the large list of files and directories, to find the ones that end in `.py`.
- ❸ And you're left with the list of unit testing scripts, because they were the only ones named `test_*.py`.

### Example 16.19. Step 3: Map filenames to module names



- ❶ As you saw in Section 4.7, "Using lambda Functions", `lambda filename: filename[:-3]` is a quick-and-dirty way of creating an inline, one-line function. This one takes a filename and returns the filename with the last three characters removed.
- ❷ `map` is a function. There's nothing magic about `map` functions as opposed to regular functions that you define with a `def` statement. You can call them with any iterable.
- ❸ Now you can apply this function to each file in the list of unit test files, using `map`.
- ❹ And the result is just what you wanted: a list of modules, as strings.

### Example 16.20. Step 4: Mapping module names to modules





```
1
```

- 1 As you saw in Section 16.6, “Dynamically importing modules”, you can use a combination of `importlib` and `sys` to map a list of module names (as seen in the previous example) to a list of module objects.
- 2 This is now a list of modules, fully accessible like any other module.
- 3 The last module in the list is the `unittest` module, just as if you had said `import unittest`.

### Example 16.21. Step 5: Loading the modules into a test suite

```
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

- 1 These are real module objects. Not only can you access them like any other module, instantiate classes and call functions, you can also use them to create test suites.
- 2 Finally, you wrap the list of module objects into one big test suite. The `unittest` module has no problem traversing this tree of nested test suites within the `unittest` module.

This introspection process is what the `unittest` module usually does for us. Remember that magic-looking `unittest` function that our individual test modules call?

### Example 16.22. Step 6: Telling `unittest` to use your test suite

```
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

- 1 Instead of letting the `unittest` module do all its magic for us, you’ve done most of it yourself. You’ve created a function (`test_suite`) that imports the modules and creates a test suite.

## 16.8. Summary

The `test_suite` program and its output should now make perfect sense.

You should now feel comfortable doing all of these things:

- Manipulating path information from the command line.
- Filtering lists using `filter` instead of list comprehensions.
- Mapping lists using `map` instead of list comprehensions.
- Dynamically importing modules.

---

Technically, the second argument to `unittest` can be any sequence, including lists, tuples, and custom classes that act like lists by defining the `__iter__` special method.

Again, I should point out that `len` can take a list, a tuple, or any object that acts like a sequence. See previous footnote about `len`

# Chapter 17. Dynamic functions

## 17.1. Diving in

I want to talk about plural nouns. Also, functions that return other functions, advanced regular expressions, and generators. Generators are new in Python 3.

If you haven't read Chapter 7, *Regular Expressions*, now would be a good time. This chapter assumes you understand the basics of regular expressions.

English is a schizophrenic language that borrows from a lot of other languages, and the rules for making singular nouns into plural nouns are a bit messy.

If you grew up in an English-speaking country or learned English in a formal school setting, you're probably familiar with the basic rules:

If a word ends in S, X, or Z, add ES. "Bass" becomes "basses", "fax" becomes "faxes", and "waltz" becomes "waltzes".<sup>1</sup>

If a word ends in a noisy H, add ES; if it ends in a silent H, just add S. What's a noisy H? One that gets combined with other letters to form a new sound.

If a word ends in Y that sounds like I, change the Y to IES; if the Y is combined with a vowel to sound like something else, just add S.

If all else fails, just add S and hope for the best.<sup>4</sup>

(I know, there are a lot of exceptions. "Man" becomes "men" and "woman" becomes "women", but "human" becomes "humans". "Mouse" becomes "mice".)

Other languages are, of course, completely different.

Let's design a module that pluralizes nouns. Start with just English nouns, and just these four rules, but keep in mind that you'll inevitably need more rules.

## 17.2. Stage 1

So you're looking at words, which at least in English are strings of characters. And you have rules that say you need to find different combinations of characters.

### Example 17.1. `pluralize.py`

```
#!/usr/bin/env python3
```

```
"""Pluralize nouns."""
```

```
import re
```

```
def pluralize(noun):
```

```
 """Return the plural form of a noun."""
```

```
 # Rule 1: Words ending in s, x, or z add es.
```

```
 # Rule 2: Words ending in a noisy h add es.
```

```
 # Rule 3: Words ending in y that sound like i change y to ies.
```

```
 # Rule 4: All other words add s.
```

```
 # Rule 5: Words ending in y that sound like e just add s.
```

```
 # Rule 6: Words ending in y that sound like i just add s.
```

❶ OK, this is a regular expression, but it uses a syntax you didn't see in Chapter 7, *Regular Expressions*. The square brackets mean "match any of these characters".

❷ This function performs regular expression-based string substitutions. Let's look at it in more detail.

### Example 17.2. Introducing `pluralize.py`

```
#!/usr/bin/env python3
```

```
"""Pluralize nouns using regular expressions. This module implements the rules from Example 17.1, but uses regular expressions to find the right rule to apply. It also uses the re.sub() function to perform the substitutions. The rules are: 1. Words ending in s, x, or z add es. 2. Words ending in a noisy h add es. 3. Words ending in y that sound like i change y to ies. 4. All other words add s. 5. Words ending in y that sound like e just add s. 6. Words ending in y that sound like i just add s. """
```

```
import re
```

```
def pluralize(noun):
```

```
 """Return the plural form of a noun using regular expressions. This function implements the rules from Example 17.1, but uses regular expressions to find the right rule to apply. It also uses the re.sub() function to perform the substitutions. The rules are: 1. Words ending in s, x, or z add es. 2. Words ending in a noisy h add es. 3. Words ending in y that sound like i change y to ies. 4. All other words add s. 5. Words ending in y that sound like e just add s. 6. Words ending in y that sound like i just add s. """
```

```
 # Rule 1: Words ending in s, x, or z add es.
```

```
 # Rule 2: Words ending in a noisy h add es.
```

```
 # Rule 3: Words ending in y that sound like i change y to ies.
```

```
 # Rule 4: All other words add s.
```

```
 # Rule 5: Words ending in y that sound like e just add s.
```

```
 # Rule 6: Words ending in y that sound like i just add s.
```



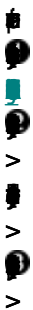
- 1 Does the string `Mountain` contain `abor`? Yes, it contains a
- 2 OK, now find `abor` and replace it with `om`. It becomes `M`
- 3 The same function turns `into` `o`
- 4 You might think this would turn `into` `o` but it doesn't. `re.sub` replaces *all* of the matches, not just the first one. So this regular expression turns

### Example 17.3. Back to `re.sub`



- 1 Back to the `re.sub` function. What are you doing? You're replacing the end of string with `o`. In other words, adding `o` to the string. You could also
- 2 Look closely, this is another new variation. The `^` as the first character inside the square brackets means something special: negation. `[^Y]`
- 3 Same pattern here: match words that end in `Y`, where the character before the `Y` is *not* `a`, `e`, `i`, or `u`. You're looking for words that end in `Y`

### Example 17.4. More on negation regular expressions



- 1 `re.sub` matches this regular expression, because it ends in `y` and `o` is not `a`, `e`, `i`, or `u`
- 2 `re.sub` does not match, because it ends in `y` and you specifically said that the character before the `y` could not be `o`. `re.sub` does not match, because it ends
- 3 `re.sub` does not match, because it does not end in `y`

### Example 17.5. More on `re.sub`





- 1 This regular expression turns `into` `and` `into` `which` is what you wanted. Note that it would also turn `into` `but` that will never happen
- 2 Just in passing, I want to point out that it is possible to combine these two regular expressions (one to find out if the rule applies, and a

Regular expression substitutions are extremely powerful, and the `\` syntax makes them even more powerful. But combining the entire operation

### 17.3. Stage 2

Now you're going to add a level of abstraction. You started by defining a list of rules: if this, then do that, otherwise go to the next rule. Let's

#### Example 17.6.



- 1 This version looks more complicated (it's certainly longer), but it does exactly the same thing: try to match four different rules, in order

- ② Using a `for` loop, you can pull out the match and apply rules two at a time (one match, one apply) from the `rules` tuple. On the first iteration of the loop, you can match the first rule, and then apply it to the `graph`.
- ③ Remember that everything in Python is an object, including functions. `rules` contains actual functions; not names of functions, but actual function objects.
- ④ On the first iteration of the `for` loop, this is equivalent to calling `rules[0].match` and so forth.

If this additional level of abstraction is confusing, try unrolling the function to see the equivalence. This `for` loop is equivalent to the following:

### Example 17.7. Unrolling the `apply_rules` function

```
def apply_rules(rules, graph):
 for rule in rules:
 match = rule.match(graph)
 if match:
 rule.apply(graph)
```

The benefit here is that that `apply_rules` function is now simplified. It takes a list of rules, defined elsewhere, and iterates through them in a generic fashion.

Now, was adding this level of abstraction worth it? Well, not yet. Let's consider what it would take to add a new rule to the function. Well, in the current code, you'd have to modify the `apply_rules` function to add the new rule.

This is really just a stepping stone to the next section. Let's move on.

## 17.4. Stage 3

Defining separate named functions for each match and apply rule isn't really necessary. You never call them directly; you define them in the `rules` tuple.

### Example 17.8. `rules` tuple

```
rules = (
 (
 (
 match_rule1,
 apply_rule1
),
 (
 match_rule2,
 apply_rule2
),
 (
 match_rule3,
 apply_rule3
)
)
)
```

```

1
2
3

```

- 1 This is the same set of rules as you defined in stage 2. The only difference is that instead of defining named functions like `land` and `by` you have `land` and `by` defined directly in the `rules` list.
- 2 Note that the `rules` function hasn't changed at all. It iterates through a set of rule functions, checks the first rule, and if it returns a true value...

Now to add a new rule, all you need to do is define the functions directly in the `rules` list itself: one match rule, and one apply rule. But defining the

## 17.5. Stage 4

Let's factor out the duplication in the code so that defining new rules can be easier.

### Example 17.9.

```

1
2
3
4

```

- 1 `rules` is a function that builds other functions dynamically. It takes `rules` and `rules` (actually it takes a tuple, but more on that in a minute), and you can use it to build other functions.
- 2 Building the apply function works the same way. The apply function is a function that takes one parameter, and calls `rules` with the `rules` and `rules` parameters.
- 3 Finally, the `rules` function returns a tuple of two values: the two functions you just created. The constants you defined within those functions are now constants within the `rules` function.

If this is incredibly confusing (and it should be, this is weird stuff), it may become clearer when you see how to use it.

### Example 17.10. Continued

```

1
2
3
4
5
6
7
8
9

```

- 1 Our pluralization rules are now defined as a series of strings (not functions). The first string is the regular expression that you would use to match a word.
- 2 This line is magic. It takes the list of strings in `rules` and turns them into a list of functions. How? By mapping the strings to the `rules` function, which

I swear I am not making this up: `rules` ends up with exactly the same list of functions as the previous example. Unroll the `rules` definition, and you'll get

### Example 17.11. Unrolling the rules definition

```

1
2
3
4
5
6

```

```
(
#
#
)
(
#
#
)
(
#
#
)
)
```

### Example 17.12. Finishing up

```
#
#
#
#
```

❶ Since the list is the same as the previous example, it should come as no surprise that the function hasn't changed. Remember, it's com

Just in case that wasn't mind-blowing enough, I must confess that there was a subtlety in the definition of that I skipped over. Let's go back

### Example 17.13. Another look at

```
#
```

❶ Notice the double parentheses? This function doesn't actually take three parameters; it actually takes one parameter, a tuple of three el

### Example 17.14. Expanding tuples when calling functions

```
#
#
#
#
#
#
#
#
#
#
#
#
```

❶ The proper way to call the function is with a tuple of three elements. When the function is called, the elements are assigned to differ

Now let's go back and see why this auto-tuple-expansion trick was necessary. was a list of tuples, and each tuple had three elements. When y

## 17.6. Stage 5

You've factored out all the duplicate code and added enough abstractions so that the pluralization rules are defined in a list of strings. The ne

First, let's create a text file that contains the rules you want. No fancy data structures, just space- (or tab-)delimited strings in three columns. T



### Example 17.15. `rules.py`

```
def pluralize(word, count):
 if count == 1:
 return word
 else:
 return word + 's'
```

Now let's see how you can use this rules file.

### Example 17.16. `pluralize.py`

```
def pluralize(word, count):
 if count == 1:
 return word
 else:
 return word + 's'
```

- 1 You're still using the closures technique here (building a function dynamically that uses variables defined outside the function), but now
- 2 Our `pluralize` function now takes an optional second parameter, `rules`, which defaults to `rules`
- 3 You use the `rules` parameter to construct a filename, then open the file and read the contents into a list. If `rules` is `None`, then you'll open the `rules.py` file, read
- 4 As you saw, each line in the file really has three values, but they're separated by whitespace (tabs or spaces, it makes no difference). Now
- 5 If `rules` is a list of tuples, then `rules` will be a list of the functions created dynamically by each call to `pluralize`. Calling `pluralize` returns a function that takes a singular
- 6 Because you're now building a combined match-and-apply function, you need to call it differently. Just call the function, and if it returns

So the improvement here is that you've completely separated the pluralization rules into an external file. Not only can the file be maintained

The downside here is that you're reading that file every time you call the `pluralize` function. I thought I could get through this entire book without using

## 17.7. Stage 6

Now you're ready to talk about generators.

### Example 17.17. `pluralize.py`

```
def pluralize(word, count):
 if count == 1:
 return word
 else:
 return word + 's'
```

5

This uses a technique called generators, which I'm not even going to try to explain until you look at a simpler example first.

### Example 17.18. Introducing generators

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

- 1 The presence of the `yield` keyword in `fib` means that this is not a normal function. It is a special kind of function which generates values one at a time.
- 2 To create an instance of the `fib` generator, just call it like any other function. Note that this does not actually execute the function code. You just get a generator object.
- 3 The `fib` function returns a generator object.
- 4 The first time you call the `next` method on the generator object, it executes the code in `fib` up to the first `yield` statement, and then returns the value `0`.
- 5 Repeatedly calling `next` on the generator object *resumes where you left off* and continues until you hit the next `yield` statement. The next line of code is `1`, so it returns `1`.
- 6 The second time you call `next` you do all the same things again, but this time `x` is now `1`. And so forth. Since `fib` sets up an infinite loop, you can call `next` as many times as you want.

### Example 17.19. Using generators instead of recursion

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

- 1 The Fibonacci sequence is a sequence of numbers where each number is the sum of the two numbers before it. It starts with `0` and `1`, goes on as `0, 1, 1, 2, 3, 5, 8, 13, 21, 34, 55, 89, 144, ...`
- 2 `fib` is the current number in the sequence, so yield it.
- 3 `fib` is the next number in the sequence, so assign that to `a` but also calculate the next value (`b + a`) and assign that to `b` for later use. Note that this is done *before* yielding `a`.

So you have a function that spits out successive Fibonacci numbers. Sure, you could do that with recursion, but this way is easier to read. Also, you can stop the sequence at any time.

### Example 17.20. Generators in loops



- 1 You can use a generator like `lines` in a `for` loop directly. The `for` loop will create the generator object and successively call the `next` method to get values.
- 2 Each time through the `for` loop, `next` gets a new value from the `yield` statement in `lines` and all you do is print it out. Once `lines` runs out of numbers (`next` gets `StopIteration`), the loop ends.

OK, let's go back to the `rules` function and see how you're using this.

### Example 17.21. Generators that generate dynamic functions



- 1 `rules` is a common idiom for reading lines from a file, one line at a time. It works because `rules` actually returns a generator whose `next` method returns the next line.
- 2 No magic here. Remember that the lines of the rules file have three values separated by whitespace, so `rules` returns a tuple of 3 values, and `yield` yields each value.
- 3 And then you yield. What do you yield? A function, built dynamically with `eval` that is actually a closure (it uses the local variables `rules` and `lines`).
- 4 Since `rules` is a generator, you can use it directly in a `for` loop. The first time through the `for` loop, you will call the `next` function, which will open the file.

What have you gained over stage 5? In stage 5, you read the entire rules file and built a list of all the possible rules before you even tried the first one.

### Further reading

- PEP 255 (<http://www.python.org/peps/pep-0255.html>) defines generators.
- Python Cookbook (<http://www.activestate.com/ASPN/Python/Cookbook/>) has many more examples of generators (<http://www.google.com/search?q=python+cookbook+generators&btnG=Search>).

## 17.8. Summary

You talked about several different advanced techniques in this chapter. Not all of them are appropriate for every situation.

You should now be comfortable with all of these techniques:

- Performing string substitution with regular expressions.
- Treating functions as objects, storing them in lists, assigning them to variables, and calling them through those variables.
- Building dynamic functions with `eval`.
- Building closures, dynamic functions that contain surrounding variables as constants.
- Building generators, resumable functions that perform incremental logic and return different values each time you call them.

Adding abstractions, building functions dynamically, building closures, and using generators can all make your code simpler, more readable,

# Chapter 18. Performance Tuning

Performance tuning is a many-splendored thing. Just because Python is an interpreted language doesn't mean you shouldn't worry about code

## 18.1. Diving in

There are so many pitfalls involved in optimizing your code, it's hard to know where to start.

Let's start here: *are you sure you need to do it at all?* Is your code really so bad? Is it worth the time to tune it? Over the lifetime of your appl

Second, *are you sure you're done coding?* Premature optimization is like spreading frosting on a half-baked cake. You spend hours or days (

This is not to say that code optimization is worthless, but you need to look at the whole system and decide whether it's the best use of your ti

Oh yes, unit tests. It should go without saying that you need a complete set of unit tests before you begin performance tuning. The last thing y

With these caveats in place, let's look at some techniques for optimizing Python code. The code in question is an implementation of the Soun

There are several subtle variations of the Soundex algorithm. This is the one used in this chapter:

Keep the first letter of the name as-is.1.

Convert the remaining letters to digits, according to a specific table:

B, F, P, and V become 1.

C, G, J, K, Q, S, X, and Z become 2.

D and T become 3.

L becomes 4.

M and N become 5.

R becomes 6.

All other letters become 9.

Remove consecutive duplicates.3.

Remove all 9s altogether.4.

If the result is shorter than four characters (the first letter plus three digits), pad the result with trailing zeros.5.

if the result is longer than four characters, discard everything after the fourth character.6.

For example, my name, **B** becomes P942695. That has no consecutive duplicates, so nothing to do there. Then you remove the 9s, leaving P426

Another example: **W** becomes W99, which becomes W9, which becomes W, which gets padded with zeros to become W000.

Here's a first attempt at a Soundex function:

### Example 18.1.

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5>).

















0 1 2 3 4 5 6 7 8 9  
 A B C D E F G H I J  
 K L M N O P Q R S  
 T U V W X Y Z  
 a b c d e f g h i j  
 k l m n o p q r s  
 t u v w x y z  
 0 1 2 3 4 5 6 7 8 9  
 A B C D E F G H I J  
 K L M N O P Q R S  
 T U V W X Y Z  
 a b c d e f g h i j  
 k l m n o p q r s  
 t u v w x y z



## Further Reading on Soundex

Soundexing and Genealogy (<http://www.avotaynu.com/soundex.html>) gives a chronology of the evolution of the Soundex and its re-

## 18.2. Using the `timeit` Module

The most important thing you need to know about optimizing Python code is that you shouldn't write your own timing function.

Timing short pieces of code is incredibly complex. How much processor time is your computer devoting to running this code? Are there things

And then there's the matter of the variations introduced by the timing framework itself. Does the Python interpreter cache method name lookups

The Python community has a saying: "Python comes with batteries included." Don't write your own timing framework. Python 2.3 comes with

### Example 18.2. Introducing `timeit`

If you have not already done so, you can download this and other examples (<http://diveintopython.org/download/diveintopython-examples-5>).



- 1 The `timeit` module defines one class, `Timer`, which takes two arguments. Both arguments are strings. The first argument is the statement you wish to time.
- 2 Once you have the `Timer` object, the easiest thing to do is call `timeit()` which calls your function 1 million times and returns the number of seconds it took.
- 3 The other major method of the `Timer` object is `repeat()` which takes two optional arguments. The first argument is the number of times to repeat the test.

You can use the `timeit` module on the command line to test an existing Python program, without modifying the code. See <http://docs.python.org/lib>

Note that `timeit()` returns a list of times. The times will almost never be identical, due to slight variations in how much processor time the Python interpreter

In fact, that's almost certainly wrong. The tests that took longer didn't take longer because of variations in your code or in the Python interpreter.

Python has a handy `min()` function that takes a list and returns the smallest value:



The `isword` module only works if you already know what piece of code you need to optimize. If you have a larger Python program and don't know

## 18.3. Optimizing Regular Expressions

The first thing the `Soundex` function checks is whether the input is a non-empty string of letters. What's the best way to do this?

If you answered "regular expressions", go sit in the corner and contemplate your bad instincts. Regular expressions are almost never the right

This code fragment from `isword` checks whether the function argument `is` a word made entirely of letters, with at least one letter (not the empty string).



How does `isword` perform? For convenience, the `test` section of the script contains this code that calls the `isword` module, sets up a timing test with three differ



So how does `isword` perform with this regular expression?



As you might expect, the algorithm takes significantly longer when called with longer names. There will be a few things we can do to narrow

The other thing to keep in mind is that we are testing a representative sample of names. `isword` is a kind of trivial case, in that it gets shorted down to

So what about that regular expression? Well, it's inefficient. Since the expression is testing for ranges of characters (`Aa` uppercase, and `a` lowercase)



`isword` is slightly faster than `isword` but nothing to get terribly excited about:



We saw in Section 15.3, “Refactoring” that regular expressions can be compiled and reused for faster results. Since this regular expression ne

```
import re
def is_valid(s):
 return re.match('^[a-zA-Z]*$', s) is not None
```

Using a compiled regular expression in `is_valid` is significantly faster:

```
import re
def is_valid(s):
 re_obj = re.compile('^[a-zA-Z]*$')
 return re_obj.match(s) is not None
```

But is this the wrong path? The logic here is simple: the input `s` needs to be non-empty, and it needs to be composed entirely of letters. Would

Here is `is_valid`

```
def is_valid(s):
 if s == '':
 return False
 for c in s:
 if not c.isalpha():
 return False
 return True
```

It turns out that this technique in `is_valid` is *not* faster than using a compiled regular expression (although it is faster than using a non-compiled regul

```
import re
def is_valid(s):
 re_obj = re.compile('^[a-zA-Z]*$')
 return re_obj.match(s) is not None
```

Why isn't `is_valid` faster? The answer lies in the interpreted nature of Python. The regular expression engine is written in C, and compiled to run nati

It turns out that Python offers an obscure string method. You can be excused for not knowing about it, since it's never been mentioned in this

This is `is_valid`

```
def is_valid(s):
 return s.isalpha()
```

How much did we gain by using this specific method in `is_valid`? Quite a bit.

```
import re
def is_valid(s):
 re_obj = re.compile('^[a-zA-Z]*$')
 return re_obj.match(s) is not None
```



### Example 18.3. Best Result So Far:





































































































## 18.4. Optimizing Dictionary Lookups

The second step of the Soundex algorithm is to convert characters to digits in a specific pattern. What's the best way to do this?

The most obvious solution is to define a dictionary with individual characters as keys and their corresponding digits as values, and do dictionary lookups for each character.



You timed it already; this is how it performs:



This code is straightforward, but is it the best solution? Calling `dict.get()` on each individual character seems inefficient; it would probably be better to use a list comprehension.

Then there's the matter of incrementally building the `string`. Incrementally building strings like this is horribly inefficient; internally, the Python is good at lists, though. It can treat a string as a list of characters automatically. And lists are easy to combine into strings again, using

Here is `str.join` which converts letters to digits by using `str.maketrans` and `str.translate`

```
def str_join(s):
 # Create a translation table
 trans = str.maketrans('abcdefghijklmnopqrstuvwxyz', '01234567890123456789')
 return s.translate(trans).join('')
```

Surprisingly, `str.join` is not faster:

```
def str_join(s):
 # Create a translation table
 trans = str.maketrans('abcdefghijklmnopqrstuvwxyz', '01234567890123456789')
 return s.translate(trans).join('')
```

The overhead of the anonymous `lambda` function kills any performance you gain by dealing with the string as a list of characters.

```
def str_join(s):
 # Create a translation table
 trans = str.maketrans('abcdefghijklmnopqrstuvwxyz', '01234567890123456789')
 return s.translate(trans).join('')
```

Using a list comprehension in `str.join` is faster than using `str.maketrans` and `str.translate` but still not faster than the original code (incrementally building a string in `str.join`)

```
def str_join(s):
 # Create a translation table
 trans = str.maketrans('abcdefghijklmnopqrstuvwxyz', '01234567890123456789')
 return s.translate(trans).join('')
```

It's time for a radically different approach. Dictionary lookups are a general purpose tool. Dictionary keys can be any length string (or many)

This is `str.join`

```
def str_join(s):
 # Create a translation table
 trans = str.maketrans('abcdefghijklmnopqrstuvwxyz', '01234567890123456789')
 return s.translate(trans).join('')
```

What the heck is going on here? `str.join` creates a translation matrix between two strings: the first argument and the second argument. In this case, the

shows that `str.join` is significantly faster than defining a dictionary and looping through the input and building the output incrementally:

```
def str_join(s):
 # Create a translation table
 trans = str.maketrans('abcdefghijklmnopqrstuvwxyz', '01234567890123456789')
 return s.translate(trans).join('')
```

You're not going to get much better than that. Python has a specialized function that does exactly what you want to do; use it and move on.

#### Example 18.4. Best Result So Far:

```
def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```


```
 def
```

```
 def
```

```
 def
```

## 18.5. Optimizing List Operations

The third step in the Soundex algorithm is eliminating consecutive duplicate digits. What's the best way to do this?

Here's the code we have so far, in 

```
def
```

```
 def
```

```
 def
```

```
 def
```

```
 def
```


Here are the performance results for 

```
def
```

```
 def
```

```
 def
```

```
 def
```

The first thing to consider is whether it's efficient to check  each time through the loop. Are list indexes expensive? Would we be better off n

To answer this question, here is `islice`:

```
islice = lambda s, i, j: iter(s[i:j])
```

`islice` does not run any faster than `l[i:j]` and may even be slightly slower (although it's not enough of a difference to say for sure):

```
islice = lambda s, i, j: iter(s[i:j])
```

Why isn't `islice` faster? It turns out that list indexes in Python are extremely efficient. Repeatedly accessing `l[i]` is no problem at all. On the other hand,

Let's try something radically different. If it's possible to treat a string as a list of characters, it should be possible to use a list comprehension

However, it is possible to create a list of index numbers using the built-in `range` function, and use those index numbers to progressively search through

Here is `islice`:

```
islice = lambda s, i, j: iter(s[i:j])
```

Is this faster? In a word, no.

```
islice = lambda s, i, j: iter(s[i:j])
```

It's possible that the techniques so far as have been "string-centric". Python can convert a string into a list of characters with a single command

Here is `uniq` which modifies a list in place to remove consecutive duplicate elements:

```
def uniq(l):
 for i in range(1, len(l)):
 if l[i] == l[i-1]:
 del l[i]
```

Is this faster than `islice` or `list`? No, in fact it's the slowest method yet:

```
def uniq(l):
```




We haven't made any progress here at all, except to try and rule out several "clever" techniques. The fastest code we've seen so far was the o

### Example 18.5. Best Result So Far:




## 18.6. Optimizing String Manipulation

The final step of the Soundex algorithm is padding short results with zeros, and truncating long results. What is the best way to do this?

This is what we have so far, taken from 



These are the results for 





The first thing to consider is replacing that regular expression with a loop. This code is from [here](#)



Is it faster? Yes it is:



But wait a minute. A loop to remove characters from a string? We can use a simple string method for that. Here's [one](#)



Is it faster? That's an interesting question. It depends on the input:



The string method in [this](#) is faster than the loop for most names, but it's actually slightly slower than [this](#) in the trivial case (of a very short name). People

Last but not least, let's examine the final two steps of the algorithm: padding short results with zeros, and truncating long results to four characters.



Why do we need a loop to pad out the result? We know in advance that we're going to truncate the result to four characters, and we know the

How much speed do we gain in [this](#) by dropping the loop? It's significant:



Finally, there is still one more thing you can do to these three lines of code to make them faster: you can combine them into one line. Take a [look](#)



Putting all this code on one line in `is` barely faster than `is`



It is also significantly less readable, and for not much performance gain. Is that worth it? I hope you have good comments. Performance isn't

## 18.7. Summary

This chapter has illustrated several important aspects of performance tuning in Python, and performance tuning in general.

- If you need to choose between regular expressions and writing a loop, choose regular expressions. The regular expression engine is
- If you need to choose between regular expressions and string methods, choose string methods. Both are compiled in C, so choose the
- General-purpose dictionary lookups are fast, but specialty functions such as `is` and string methods such as `is` are faster. If Python has a
- Don't be too clever. Sometimes the most obvious algorithm is also the fastest.
- Don't sweat it too much. Performance isn't everything.

I can't emphasize that last point strongly enough. Over the course of this chapter, you made this function three times faster and saved 20 seconds



# Appendix A. Further reading

## Chapter 1. Installing Python

## Chapter 2. Your First Python Program

### 2.3. Documenting Functions

PEP 257 (<http://www.python.org/peps/pep-0257.html>) defines [conventions](#).

*Python Style Guide* (<http://www.python.org/doc/essays/styleguide.html>) discusses how to write a good [style](#).

*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) discusses conventions for spacing in [code](#) (<http://www.python.org/doc/essays/styleguide.html>).

### 2.4.2. What's an Object?

*Python Reference Manual* (<http://www.python.org/doc/current/ref/>) explains exactly what it means to say that everything in Python is an object.

*eff-bot* (<http://www.effbot.org/guides/>) summarizes Python objects (<http://www.effbot.org/guides/python-objects.htm>).

### 2.5. Indenting Code

*Python Reference Manual* (<http://www.python.org/doc/current/ref/>) discusses cross-platform indentation issues and shows examples of good indentation style.

*Python Style Guide* (<http://www.python.org/doc/essays/styleguide.html>) discusses good indentation style.

### 2.6. Testing Modules

*Python Reference Manual* (<http://www.python.org/doc/current/ref/>) discusses the low-level details of importing modules ([http://www.python.org/doc/current/ref/import.html](#)).

## Chapter 3. Native Datatypes

### 3.1.3. Deleting Items From Dictionaries

*How to Think Like a Computer Scientist* (<http://www.ibiblio.org/obp/thinkCSPy/>) teaches about dictionaries and shows how to delete items from a dictionary.

Python Knowledge Base (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>) has a lot of example code using dictionaries.

Python Cookbook (<http://www.activestate.com/ASPN/Python/Cookbook/>) discusses how to sort the values of a dictionary.

*Python Library Reference* (<http://www.python.org/doc/current/lib/>) summarizes all the dictionary methods (<http://www.python.org/doc/current/lib/dict-objects.html>).

### 3.2.5. Using List Operators

*How to Think Like a Computer Scientist* (<http://www.ibiblio.org/obp/thinkCSPy/>) teaches about lists and makes an important distinction between lists and arrays.

*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) shows how to use lists as stacks and queues (<http://www.python.org/doc/current/tut/tut.html>).

Python Knowledge Base (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>) answers common questions about lists.

*Python Library Reference* (<http://www.python.org/doc/current/lib/>) summarizes all the list methods (<http://www.python.org/doc/current/lib/list-objects.html>).

### 3.3. Introducing Tuples

*How to Think Like a Computer Scientist* (<http://www.ibiblio.org/obp/thinkCSPy/>) teaches about tuples and shows how to compare tuples.

Python Knowledge Base (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>) shows how to sort a tuple (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>).

*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) shows how to define a tuple with one element (<http://www.python.org/doc/current/tut/tut.html>).

### 3.4.2. Assigning Multiple Values at Once

*Python Reference Manual* (<http://www.python.org/doc/current/ref/>) shows examples of when you can skip the line continuation character.

*How to Think Like a Computer Scientist* (<http://www.ibiblio.org/obp/thinkCSPy/>) shows how to use multi-variable assignment.

### 3.5. Formatting Strings

*Python Library Reference* (<http://www.python.org/doc/current/lib/>) summarizes all the string formatting format characters.

*Effective AWK Programming* ([http://www.gnu.org.org:8080/cgi-bin/info2www?\(gawk\)Top](http://www.gnu.org.org:8080/cgi-bin/info2www?(gawk)Top)) discusses all the format characters.

### 3.6. Mapping Lists

*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) discusses another way to map lists using the built-in `map` function.

*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) shows how to do nested list comprehensions (<http://www.python.org/doc/current/tut/tut.html>).

### 3.7. Joining Lists and Splitting Strings

Python Knowledge Base (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>) answers common questions about strings.

*Python Library Reference* (<http://www.python.org/doc/current/lib/>) summarizes all the string methods (<http://www.python.org/doc/current/lib/string-objects.html>).

*Python Library Reference* (<http://www.python.org/doc/current/lib/>) documents the `unicodedata` module (<http://www.python.org/doc/current/lib/unicodedata-module.html>).

*The Whole Python FAQ* (<http://www.python.org/doc/FAQ.html>) explains why `is` is a string method (<http://www.python.org/doc/faq/string.html>).

## Chapter 4. The Power Of Introspection

### 4.2. Using Optional and Named Arguments

- Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) discusses exactly when and how default arguments are evaluated.
- 4.3.3. Built-In Functions  
*Python Library Reference* (<http://www.python.org/doc/current/lib/>) documents all the built-in functions (<http://www.python.org/doc/current/lib/#built-in-functions>).
- 4.5. Filtering Lists  
*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) discusses another way to filter lists using the built-in `filter` function (<http://www.python.org/doc/current/tut/tut.html#filter>).
- 4.6.1. Using the and-or Trick  
 Python Cookbook (<http://www.activestate.com/ASPN/Python/Cookbook/>) discusses alternatives to the `and-or` trick (<http://www.activestate.com/ASPN/Python/Cookbook/recipe/10674.html>).
- 4.7.1. Real-World lambda Functions  
 Python Knowledge Base (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>) discusses using `lambda` to call functions in `lambda` expressions (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>).
- Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) shows how to access outside variables from inside a `lambda` function (<http://www.python.org/doc/current/tut/tut.html#lambda>).
- The Whole Python FAQ* (<http://www.python.org/doc/FAQ.html>) has examples of obfuscated one-liners using `lambda` (<http://www.python.org/doc/FAQ.html#lambda>).

## Chapter 5. Objects and Object-Orientation

- 5.2. Importing Modules Using `from module import`  
 eff-bot (<http://www.effbot.org/guides/>) has more to say on `from module import` (<http://www.effbot.org/guides/import-confusion.htm>).
- Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) discusses advanced import techniques, including `from module import *` (<http://www.python.org/doc/current/tut/tut.html#from-module-import>).
- 5.3.2. Knowing When to Use `self` and `__init__`  
*Learning to Program* (<http://www.freenetpages.co.uk/hp/alan.gauld/>) has a gentler introduction to classes (<http://www.freenetpages.co.uk/hp/alan.gauld/learning-to-program.html>).
- How to Think Like a Computer Scientist* (<http://www.ibiblio.org/obp/thinkCSpy/>) shows how to use classes to model computation (<http://www.ibiblio.org/obp/thinkCSpy/>).
- Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) has an in-depth look at classes, namespaces, and inheritance (<http://www.python.org/doc/current/tut/tut.html#classes>).
- Python Knowledge Base (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>) answers common questions about classes (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>).
- 5.4.1. Garbage Collection  
*Python Library Reference* (<http://www.python.org/doc/current/lib/>) summarizes built-in attributes like `__del__` (<http://www.python.org/doc/current/lib/#del>).
- Python Library Reference* (<http://www.python.org/doc/current/lib/>) documents the `gc` module (<http://www.python.org/doc/current/lib/#gc-module>).
- 5.5. Exploring `UserDict`: A Wrapper Class  
*Python Library Reference* (<http://www.python.org/doc/current/lib/>) documents the `UserDict` module (<http://www.python.org/doc/current/lib/#userdict-module>).
- 5.7. Advanced Special Class Methods  
*Python Reference Manual* (<http://www.python.org/doc/current/ref/>) documents all the special class methods (<http://www.python.org/doc/current/ref/>).
- 5.9. Private Functions  
*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) discusses the inner workings of private variables (<http://www.python.org/doc/current/tut/tut.html#private>).

## Chapter 6. Exceptions and File Handling

- 6.1.1. Using Exceptions For Other Purposes  
*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) discusses defining and raising your own exceptions, and `try` (<http://www.python.org/doc/current/tut/tut.html#try>).
- Python Library Reference* (<http://www.python.org/doc/current/lib/>) summarizes all the built-in exceptions (<http://www.python.org/doc/current/lib/#built-in-exceptions>).
- Python Library Reference* (<http://www.python.org/doc/current/lib/>) documents the `getpass` module (<http://www.python.org/doc/current/lib/#getpass-module>).
- Python Library Reference* (<http://www.python.org/doc/current/lib/>) documents the `random` module (<http://www.python.org/doc/current/lib/#random-module>).
- Python Reference Manual* (<http://www.python.org/doc/current/ref/>) discusses the inner workings of the `try` block (<http://www.python.org/doc/current/ref/#try>).
- 6.2.4. Writing to Files  
*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) discusses reading and writing files, including how to read from a file (<http://www.python.org/doc/current/tut/tut.html#reading>).
- eff-bot (<http://www.effbot.org/guides/>) discusses efficiency and performance of various ways of reading a file (<http://www.effbot.org/guides/reading-files.htm>).
- Python Knowledge Base (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>) answers common questions about file objects (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>).
- Python Library Reference* (<http://www.python.org/doc/current/lib/>) summarizes all the file object methods (<http://www.python.org/doc/current/lib/#file-object-methods>).
- 6.4. Using `sys.modules`  
*Python Tutorial* (<http://www.python.org/doc/current/tut/tut.html>) discusses exactly when and how default arguments are evaluated (<http://www.python.org/doc/current/tut/tut.html#default-arguments>).
- Python Library Reference* (<http://www.python.org/doc/current/lib/>) documents the `sys` module (<http://www.python.org/doc/current/lib/#sys-module>).
- 6.5. Working with Directories  
 Python Knowledge Base (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>) answers questions about the `os` module (<http://www.faqs.com/knowledge-base/index.phtml/fid/199/>).
- Python Library Reference* (<http://www.python.org/doc/current/lib/>) documents the `os` module (<http://www.python.org/doc/current/lib/#os-module>).

## Chapter 7. Regular Expressions

## 7.6. Case study: Parsing Phone Numbers

Regular Expression HOWTO (<http://py-howto.sourceforge.net/regex/regex.html>) teaches about regular expressions and how to use them.  
*Python Library Reference* (<http://www.python.org/doc/current/lib/>) summarizes the `re` module (<http://www.python.org/doc/current/lib/module-re.html>).

## Chapter 8. HTML Processing

### 8.4. Introducing BaseHTMLProcessor.py

W3C (<http://www.w3.org/>) discusses character and entity references (<http://www.w3.org/TR/REC-html40/charset.html#ent>).  
*Python Library Reference* (<http://www.python.org/doc/current/lib/>) confirms your suspicions that the `urllib` module (<http://www.python.org/doc/current/lib/module-urllib.html>).

### 8.9. Putting it all together

You thought I was kidding about the server-side scripting idea. So did I, until I found this web-based dialectizer (<http://rink.com/dialectizer/>).

## Chapter 9. XML Processing

### 9.4. Unicode

Unicode.org (<http://www.unicode.org/>) is the home page of the unicode standard, including a brief technical introduction (Unicode Tutorial ([http://www.reportlab.com/i18n/python\\_unicode\\_tutorial.html](http://www.reportlab.com/i18n/python_unicode_tutorial.html)) has some more examples of how to use Python's unicode support.  
PEP 263 (<http://www.python.org/peps/pep-0263.html>) goes into more detail about how and when to define a character encoding.

## Chapter 10. Scripts and Streams

## Chapter 11. HTTP Web Services

### 11.1. Diving in

Paul Prescod believes that pure HTTP web services are the future of the Internet (<http://webservicex.com/pub/a/ws/2000/06/06/>).

## Chapter 12. SOAP Web Services

### 12.1. Diving In

<http://www.xmethods.net/> is a repository of public access SOAP web services.

The SOAP specification (<http://www.w3.org/TR/soap/>) is surprisingly readable, if you like that sort of thing.

### 12.8. Troubleshooting SOAP Web Services

New developments for SOAPpy (<http://www-106.ibm.com/developerworks/webservices/library/ws-pyth17.html>) steps through some common problems.

## Chapter 13. Unit Testing

### 13.1. Introduction to Roman numerals

This site (<http://www.wilkiecollins.demon.co.uk/roman/front.htm>) has more on Roman numerals, including a fascinating history.

### 13.3. Introducing romantest.py

The PyUnit home page (<http://pyunit.sourceforge.net/>) has an in-depth discussion of using the `PyUnit` framework (<http://pyunit.sourceforge.net/>).  
The PyUnit FAQ (<http://pyunit.sourceforge.net/pyunit.html>) explains why test cases are stored separately (<http://pyunit.sourceforge.net/pyunit.html#faq>).  
*Python Library Reference* (<http://www.python.org/doc/current/lib/>) summarizes the `unittest` module (<http://www.python.org/doc/current/lib/module-unittest.html>).  
ExtremeProgramming.org (<http://www.extremeprogramming.org/>) discusses why you should write unit tests (<http://www.extremeprogramming.org/why-unit-tests.html>).  
The Portland Pattern Repository (<http://www.c2.com/cgi/wiki/>) has an ongoing discussion of unit tests (<http://www.c2.com/cgi/wiki/Unit+Testing>).

## Chapter 14. Test-First Programming

## Chapter 15. Refactoring

### 15.5. Summary

XProgramming.com (<http://www.xprogramming.com/>) has links to download unit testing frameworks (<http://www.xprogramming.com/unit-testing/>).

## Chapter 16. Functional Programming

## Chapter 17. Dynamic functions

### 17.7. plural.py, stage 6

PEP 255 (<http://www.python.org/peps/pep-0255.html>) defines generators.

Python Cookbook (<http://www.activestate.com/ASPN/Python/Cookbook/>) has many more examples of generators (<http://w>

## Chapter 18. Performance Tuning

### 18.1. Diving in

Soundexing and Genealogy (<http://www.avotaynu.com/soundex.html>) gives a chronology of the evolution of the Soundex

# Appendix B. A 5-minute review

## Chapter 1. Installing Python

### 1.1. Which Python is right for you?

The first thing you need to do with Python is install it. Or do you?

### 1.2. Python on Windows

On Windows, you have a couple choices for installing Python.

### 1.3. Python on Mac OS X

On Mac OS X, you have two choices for installing Python: install it, or don't install it. You probably want to install it.

### 1.4. Python on Mac OS 9

Mac OS 9 does not come with any version of Python, but installation is very simple, and there is only one choice.

### 1.5. Python on RedHat Linux

Download the latest Python RPM by going to <http://www.python.org/ftp/python/> and selecting the highest version number

### 1.6. Python on Debian GNU/Linux

If you are lucky enough to be running Debian GNU/Linux, you install Python through the **apt** command.

### 1.7. Python Installation from Source

If you prefer to build from source, you can download the Python source code from <http://www.python.org/ftp/python/>. Sele

### 1.8. The Interactive Shell

Now that you have Python installed, what's this interactive shell thing you're running?

### 1.9. Summary

You should now have a version of Python installed that works for you.

## Chapter 2. Your First Python Program

### 2.1. Diving in

Here is a complete, working Python program.

### 2.2. Declaring Functions

Python has functions like most other languages, but it does not have separate header files like C++ or sections like Pascal.

### 2.3. Documenting Functions

You can document a Python function by giving it a 

## 2.4. Everything Is an Object

A function, like everything else in Python, is an object.

## 2.5. Indenting Code

Python functions have no explicit `begin` and no curly braces to mark where the function code starts and stops. The only delin

## 2.6. Testing Modules

Python modules are objects and have several useful attributes. You can use this to easily test your modules as you write the

# Chapter 3. Native Datatypes

## 3.1. Introducing Dictionaries

One of Python's built-in datatypes is the dictionary, which defines one-to-one relationships between keys and values.

## 3.2. Introducing Lists

Lists are Python's workhorse datatype. If your only experience with lists is arrays in Visual Basic or (God forbid) the datas

## 3.3. Introducing Tuples

A tuple is an immutable list. A tuple can not be changed in any way once it is created.

## 3.4. Declaring variables

Python has local and global variables like most other languages, but it has no explicit variable declarations. Variables sprin

## 3.5. Formatting Strings

Python supports formatting values into strings. Although this can include very complicated expressions, the most basic usa

## 3.6. Mapping Lists

One of the most powerful features of Python is the list comprehension, which provides a compact way of mapping a list int

## 3.7. Joining Lists and Splitting Strings

You have a list of key-value pairs in the form `key=value` and you want to join them into a single string. To join any list of strings into

## 3.8. Summary

The `program` and its output should now make perfect sense.

# Chapter 4. The Power Of Introspection

## 4.1. Diving In

Here is a complete, working Python program. You should understand a good deal about it just by looking at it. The number

## 4.2. Using Optional and Named Arguments

Python allows function arguments to have default values; if the function is called without the argument, the argument gets

#### 4.3. Using type, str, dir, and Other Built-In Functions

Python has a small set of extremely useful built-in functions. All other functions are partitioned off into modules. This was

#### 4.4. Getting Object References With getattr

You already know that Python functions are objects. What you don't know is that you can get a reference to a function with

#### 4.5. Filtering Lists

As you know, Python has powerful capabilities for mapping lists into other lists, via list comprehensions (Section 3.6, "Ma

#### 4.6. The Peculiar Nature of and and or

In Python, `and` and `or` perform boolean logic as you would expect, but they do not return boolean values; instead, they return one

#### 4.7. Using lambda Functions

Python supports an interesting syntax that lets you define one-line mini-functions on the fly. Borrowed from Lisp, these so-

#### 4.8. Putting It All Together

The last line of code, the only one you haven't deconstructed yet, is the one that does all the work. But by now the work is

#### 4.9. Summary

The program and its output should now make perfect sense.

### Chapter 5. Objects and Object-Orientation

#### 5.1. Diving In

Here is a complete, working Python program. Read the of the module, the classes, and the functions to get an overview of

#### 5.2. Importing Modules Using from module import

Python has two ways of importing modules. Both are useful, and you should know when to use each. One way, you've alr

#### 5.3. Defining Classes

Python is fully object-oriented: you can define your own classes, inherit from your own or built-in classes, and instantiate t

#### 5.4. Instantiating Classes

Instantiating classes in Python is straightforward. To instantiate a class, simply call the class as if it were a function, passing

#### 5.5. Exploring UserDict: A Wrapper Class

As you've seen, is a class that acts like a dictionary. To explore this further, let's look at the class in the module, which is

#### 5.6. Special Class Methods

In addition to normal class methods, there are a number of special methods that Python classes can define. Instead of being

### 5.7. Advanced Special Class Methods

Python has more special methods than just `__init__` and `__str__`. Some of them let you emulate functionality that you may not even know a

### 5.8. Introducing Class Attributes

You already know about data attributes, which are variables owned by a specific instance of a class. Python also supports c

### 5.9. Private Functions

Unlike in most languages, whether a Python function, method, or attribute is private or public is determined entirely by its

### 5.10. Summary

That's it for the hard-core object trickery. You'll see a real-world application of special class methods in Chapter 12, which

## Chapter 6. Exceptions and File Handling

### 6.1. Handling Exceptions

Like many other programming languages, Python has exception handling via `try` blocks.

### 6.2. Working with File Objects

Python has a built-in function, `open`, for opening a file on disk. `open` returns a file object, which has methods and attributes for getting

### 6.3. Iterating with `for` Loops

Like most other languages, Python has `for` loops. The only reason you haven't seen them until now is that Python is good at so

### 6.4. Using `sys.modules`

Modules, like everything else in Python, are objects. Once imported, you can always get a reference to a module through th

### 6.5. Working with Directories

The `os` module has several functions for manipulating files and directories. Here, we're looking at handling pathnames and lis

### 6.6. Putting It All Together

Once again, all the dominoes are in place. You've seen how each line of code works. Now let's step back and see how it al

### 6.7. Summary

The `try` program introduced in Chapter 5 should now make perfect sense.

## Chapter 7. Regular Expressions

### 7.1. Diving In

If what you're trying to do can be accomplished with string functions, you should use them. They're fast and simple and ea



## 7.2. Case Study: Street Addresses

This series of examples was inspired by a real-life problem I had in my day job several years ago, when I needed to scrub a

## 7.3. Case Study: Roman Numerals

You've most likely seen Roman numerals, even if you didn't recognize them. You may have seen them in copyrights of old

## 7.4. Using the {n,m} Syntax

In the previous section, you were dealing with a pattern where the same character could be repeated up to three times. Ther

## 7.5. Verbose Regular Expressions

So far you've just been dealing with what I'll call "compact" regular expressions. As you've seen, they are difficult to read

## 7.6. Case study: Parsing Phone Numbers

So far you've concentrated on matching whole patterns. Either the pattern matches, or it doesn't. But regular expressions ar

## 7.7. Summary

This is just the tiniest tip of the iceberg of what regular expressions can do. In other words, even though you're completely

# Chapter 8. HTML Processing

## 8.1. Diving in

I often see questions on comp.lang.python (<http://groups.google.com/groups?group=comp.lang.python>) like "How can I lis

## 8.2. Introducing sgmlib.py

HTML processing is broken into three steps: breaking down the HTML into its constituent pieces, fiddling with the pieces,

## 8.3. Extracting data from HTML documents

To extract data from HTML documents, subclass the `HTMLParser` class and define methods for each tag or entity you want to capture.

## 8.4. Introducing BaseHTMLProcessor.py

`BaseHTMLProcessor` doesn't produce anything by itself. It parses and parses and parses, and it calls a method for each interesting thing it finds,

## 8.5. locals and globals

Let's digress from HTML processing for a minute and talk about how Python handles variables. Python has two built-in fu

## 8.6. Dictionary-based string formatting

There is an alternative form of string formatting that uses dictionaries instead of tuples of values.

## 8.7. Quoting attribute values

This is generally a good idea. It's a good idea to quote attribute values in HTML. It's a good idea to quote attribute values in HTML. It's a good idea to quote attribute values in HTML.

## 8.8. Introducing dialect.py

`Dialect` is a simple (and silly) descendant of `HTMLParser`. It runs blocks of text through a series of substitutions, but it makes sure that anything

## 8.9. Putting it all together

It's time to put everything you've learned so far to good use. I hope you were paying attention.

## 8.10. Summary

Python provides you with a powerful tool, `lxml`, to manipulate HTML by turning its structure into an object model. You can use

## Chapter 9. XML Processing

### 9.1. Diving in

There are two basic ways to work with XML. One is called SAX ("Simple API for XML"), and it works by reading the XML

### 9.2. Packages

Actually parsing an XML document is very simple: one line of code. However, before you get to that line of code, you need

### 9.3. Parsing XML

As I was saying, actually parsing an XML document is very simple: one line of code. Where you go from there is up to you

### 9.4. Unicode

Unicode is a system to represent characters from all the world's different languages. When Python parses an XML document

### 9.5. Searching for elements

Traversing XML documents by stepping through each node can be tedious. If you're looking for something in particular, b

### 9.6. Accessing element attributes

XML elements can have one or more attributes, and it is incredibly simple to access them once you have parsed an XML d

### 9.7. Segue

OK, that's it for the hard-core XML stuff. The next chapter will continue to use these same example programs, but focus on

## Chapter 10. Scripts and Streams

### 10.1. Abstracting input sources

One of Python's greatest strengths is its dynamic binding, and one powerful use of dynamic binding is the *file-like object*.

### 10.2. Standard input, output, and error

UNIX users are already familiar with the concept of standard input, standard output, and standard error. This section is for

### 10.3. Caching node lookups

`lxml` employs several tricks which may or may not be useful to you in your XML processing. The first one takes advantage of th

### 10.4. Finding direct children of a node

Another useful technique when parsing XML documents is finding all the direct child elements of a particular element. For i

## 10.5. Creating separate handlers by node type

The third useful XML processing tip involves separating your code into logical functions, based on node types and element

## 10.6. Handling command-line arguments

Python fully supports creating programs that can be run on the command line, complete with command-line arguments and

## 10.7. Putting it all together

You've covered a lot of ground. Let's step back and see how all the pieces fit together.

## 10.8. Summary

Python comes with powerful libraries for parsing and manipulating XML documents. The `lxml` module takes an XML file and parses it

# Chapter 11. HTTP Web Services

## 11.1. Diving in

You've learned about HTML processing and XML processing, and along the way you saw how to download a web page and

## 11.2. How not to fetch data over HTTP

Let's say you want to download a resource over HTTP, such as a syndicated Atom feed. But you don't just want to download

## 11.3. Features of HTTP

There are five important features of HTTP which you should support.

## 11.4. Debugging HTTP web services

First, let's turn on the debugging features of Python's HTTP library and see what's being sent over the wire. This will be useful

## 11.5. Setting the User-Agent

The first step to improving your HTTP web services client is to identify yourself properly with a `User-Agent`. To do that, you need to modify

## 11.6. Handling Last-Modified and ETag

Now that you know how to add custom HTTP headers to your web service requests, let's look at adding support for `Last-Modified` and `ETag`.

## 11.7. Handling redirects

You can support permanent and temporary redirects using a different kind of custom URL handler.

## 11.8. Handling compressed data

The last important HTTP feature you want to support is compression. Many web services have the ability to send data compressed

## 11.9. Putting it all together

You've seen all the pieces for building an intelligent HTTP web services client. Now let's see how they all fit together.

## 11.10. Summary

The `url` and its functions should now make perfect sense.

## Chapter 12. SOAP Web Services

### 12.1. Diving In

You use Google, right? It's a popular search engine. Have you ever wished you could programmatically access Google search?

### 12.2. Installing the SOAP Libraries

Unlike the other code in this book, this chapter relies on libraries that do not come pre-installed with Python.

### 12.3. First Steps with SOAP

The heart of SOAP is the ability to call remote functions. There are a number of public access SOAP servers that provide services.

### 12.4. Debugging SOAP Web Services

The SOAP libraries provide an easy way to see what's going on behind the scenes.

### 12.5. Introducing WSDL

The `Proxy` class proxies local method calls and transparently turns them into invocations of remote SOAP methods. As you've seen, this is done by using the `__getattr__` method.

### 12.6. Introspecting SOAP Web Services with WSDL

Like many things in the web services arena, WSDL has a long and checkered history, full of political strife and intrigue. I won't go into the details here.

### 12.7. Searching Google

Let's finally turn to the sample code that you saw at the beginning of this chapter, which does something more useful and interesting than just searching Google.

### 12.8. Troubleshooting SOAP Web Services

Of course, the world of SOAP web services is not all happiness and light. Sometimes things go wrong.

### 12.9. Summary

SOAP web services are very complicated. The specification is very ambitious and tries to cover many different use cases for web services.

## Chapter 13. Unit Testing

### 13.1. Introduction to Roman numerals

In previous chapters, you "dived in" by immediately looking at code and trying to understand it as quickly as possible. Now we'll take a different approach.

### 13.2. Diving in

Now that you've completely defined the behavior you expect from your conversion functions, you're going to do something more interesting.

### 13.3. Introducing `romantest.py`

This is the complete test suite for your Roman numeral conversion functions, which are yet to be written but will eventually

#### 13.4. Testing for success

The most fundamental part of unit testing is constructing individual test cases. A test case answers a single question about t

#### 13.5. Testing for failure

It is not enough to test that functions succeed when given good input; you must also test that they fail when given bad input

#### 13.6. Testing for sanity

Often, you will find that a unit of code contains a set of reciprocal functions, usually in the form of conversion functions w

### Chapter 14. Test-First Programming

#### 14.1. roman.py, stage 1

Now that the unit tests are complete, it's time to start writing the code that the test cases are attempting to test. You're goin

#### 14.2. roman.py, stage 2

Now that you have the framework of the module laid out, it's time to start writing code and passing test cases.

#### 14.3. roman.py, stage 3

Now that behaves correctly with good input (integers from 1 to 9) it's time to make it behave correctly with bad input (every

#### 14.4. roman.py, stage 4

Now that is done, it's time to start coding. Thanks to the rich data structure that maps individual Roman numerals to integ

#### 14.5. roman.py, stage 5

Now that works properly with good input, it's time to fit in the last piece of the puzzle: making it work properly with bad i

### Chapter 15. Refactoring

#### 15.1. Handling bugs

Despite your best efforts to write comprehensive unit tests, bugs happen. What do I mean by "bug"? A bug is a test case yo

#### 15.2. Handling changing requirements

Despite your best efforts to pin your customers to the ground and extract exact requirements from them on pain of horrible

#### 15.3. Refactoring

The best thing about comprehensive unit testing is not the feeling you get when all your test cases finally pass, or even the

#### 15.4. Postscript

A clever reader read the previous section and took it to the next level. The biggest headache (and performance drain) in the

## 15.5. Summary

Unit testing is a powerful concept which, if properly implemented, can both reduce maintenance costs and increase flexibility.

## Chapter 16. Functional Programming

### 16.1. Diving in

In Chapter 13, *Unit Testing*, you learned about the philosophy of unit testing. In Chapter 14, *Test-First Programming*, you

### 16.2. Finding the path

When running Python scripts from the command line, it is sometimes useful to know where the currently running script is located.

### 16.3. Filtering lists revisited

You're already familiar with using list comprehensions to filter lists. There is another way to accomplish this same thing, with `filter()`.

### 16.4. Mapping lists revisited

You're already familiar with using list comprehensions to map one list into another. There is another way to accomplish this, with `map()`.

### 16.5. Data-centric programming

By now you're probably scratching your head wondering why this is better than using `for` loops and straight function calls. And the answer is, it's not.

### 16.6. Dynamically importing modules

OK, enough philosophizing. Let's talk about dynamically importing modules.

### 16.7. Putting it all together

You've learned enough now to deconstruct the first seven lines of this chapter's code sample: reading a directory and importing modules.

### 16.8. Summary

The `plural.py` program and its output should now make perfect sense.

## Chapter 17. Dynamic functions

### 17.1. Diving in

I want to talk about plural nouns. Also, functions that return other functions, advanced regular expressions, and generators.

### 17.2. plural.py, stage 1

So you're looking at words, which at least in English are strings of characters. And you have rules that say you need to find the plural of a word.

### 17.3. plural.py, stage 2

Now you're going to add a level of abstraction. You started by defining a list of rules: if this, then do that, otherwise go to the next rule.

### 17.4. plural.py, stage 3

Defining separate named functions for each match and apply rule isn't really necessary. You never call them directly; you call `pluralize`.

#### 17.5. plural.py, stage 4

Let's factor out the duplication in the code so that defining new rules can be easier.

#### 17.6. plural.py, stage 5

You've factored out all the duplicate code and added enough abstractions so that the pluralization rules are defined in a list.

#### 17.7. plural.py, stage 6

Now you're ready to talk about generators.

#### 17.8. Summary

You talked about several different advanced techniques in this chapter. Not all of them are appropriate for every situation.

### Chapter 18. Performance Tuning

#### 18.1. Diving in

There are so many pitfalls involved in optimizing your code, it's hard to know where to start.

#### 18.2. Using the `timeit` Module

The most important thing you need to know about optimizing Python code is that you shouldn't write your own timing function.

#### 18.3. Optimizing Regular Expressions

The first thing the `Soundex` function checks is whether the input is a non-empty string of letters. What's the best way to do this?

#### 18.4. Optimizing Dictionary Lookups

The second step of the `Soundex` algorithm is to convert characters to digits in a specific pattern. What's the best way to do this?

#### 18.5. Optimizing List Operations

The third step in the `Soundex` algorithm is eliminating consecutive duplicate digits. What's the best way to do this?

#### 18.6. Optimizing String Manipulation

The final step of the `Soundex` algorithm is padding short results with zeros, and truncating long results. What is the best way to do this?

#### 18.7. Summary

This chapter has illustrated several important aspects of performance tuning in Python, and performance tuning in general.

# Appendix C. Tips and tricks

## Chapter 1. Installing Python

## Chapter 2. Your First Python Program

### 2.1. Diving in

In the ActivePython IDE on Windows, you can run the Python program you're editing by choosing File->Run... (**Ctrl-R**). Output is

In the Python IDE on Mac OS, you can run a Python program with Python->Run window... (**Cmd-R**), but there is an important option

On UNIX-compatible systems (including Mac OS X), you can run a Python program from the command line:

### 2.2. Declaring Functions

In Visual Basic, functions (that return a value) start with **Function** and subroutines (that do not return a value) start with **Sub**. There are no subroutines

In Java, C++, and other statically-typed languages, you must specify the datatype of the function return value and each function argument

### 2.3. Documenting Functions

Triple quotes are also an easy way to define a string with both single and double quotes, like in Perl.

Many Python IDEs use the docstring to provide context-sensitive documentation, so that when you type a function name, its documentation appears as a tooltip

### 2.4. Everything Is an Object

In Python is like in Perl. Once you import a Python module, you access its functions with `module.function`. Once you import a Perl module, you access its functions with `$module->function`

### 2.5. Indenting Code

Python uses carriage returns to separate statements and a colon and indentation to separate code blocks. C++ and Java use semicolons to separate statements

### 2.6. Testing Modules

Like C, Python uses `==` for comparison and `=` for assignment. Unlike C, Python does not support in-line assignment, so there's no chance of

On MacPython, there is an additional step to make the trick work. Pop up the module's options menu by clicking the black triangle next to the module name

## Chapter 3. Native Datatypes

### 3.1. Introducing Dictionaries

A dictionary in Python is like a hash in Perl. In Perl, variables that store hashes always start with a `%` character. In Python, variables containing hashes



A dictionary in Python is like an instance of the `dict` class in Java.

A dictionary in Python is like an instance of the `Object` in Visual Basic.

### 3.1.2. Modifying Dictionaries

Dictionaries have no concept of order among elements. It is incorrect to say that the elements are "out of order"; they are simply unordered.

### 3.2. Introducing Lists

A list in Python is like an array in Perl. In Perl, variables that store arrays always start with the `@` character; in Python, variables can be any name.

A list in Python is much more than an array in Java (although it can be used as one if that's really all you want out of life). A better analogy is a C++ `vector`.

### 3.2.3. Searching Lists

Before version 2.2.1, Python had no separate boolean datatype. To compensate for this, Python accepted almost anything in a boolean context:

- 0 is false; all other numbers are true.

- An empty string (`''`) is false, all other strings are true.

- An empty list (`[]`) is false; all other lists are true.

- An empty tuple (`()`) is false; all other tuples are true.

- An empty dictionary (`{}`) is false; all other dictionaries are true.

These rules still apply in Python 2.2.1 and beyond, but now you can also use an actual boolean, which has a value of `True` or `False`. Note the capitalization.

### 3.3. Introducing Tuples

Tuples can be converted into lists, and vice-versa. The built-in `tuple()` function takes a list and returns a tuple with the same elements, and the built-in `list()` function takes a tuple and returns a list with the same elements.

### 3.4. Declaring variables

When a command is split among several lines with the line-continuation marker (`"\"`), the continued lines can be indented in any manner.

### 3.5. Formatting Strings

String formatting in Python uses the same syntax as the `printf` function in C.

### 3.7. Joining Lists and Splitting Strings

`join()` works only on lists of strings; `split()` does not do any type coercion. Joining a list that has one or more non-string elements will raise an error.

`find()` is a useful technique when you want to search a string for a substring and then work with everything before the substring (which ends with the substring).

## Chapter 4. The Power Of Introspection

### 4.2. Using Optional and Named Arguments

The only thing you need to do to call a function is specify a value (somehow) for each required argument; the manner and order in which you specify the arguments is up to you.

### 4.3.3. Built-In Functions

Python comes with excellent reference manuals, which you should peruse thoroughly to learn all the modules Python has to offer. Be sure to read the `Python Library Reference`.

### 4.7. Using lambda Functions

Functions are a matter of style. Using them is never required; anywhere you could use them, you could define a separate normal function.

#### 4.8. Putting It All Together

In SQL, you must use `IS NULL` to compare a null value. In Python, you can use either `is` or `is not` but `is` is faster.

## Chapter 5. Objects and Object-Oriented

### 5.2. Importing Modules Using `from module import`

In Python is like in Perl; in Python is like in Perl.

In Python is like in Java; in Python is like in Java.

Use sparingly, because it makes it difficult to determine where a particular function or attribute came from, and that makes debugging harder.

### 5.3. Defining Classes

The `class` statement in Python is like an empty set of braces (`{ }` in Java or C).

In Python, the ancestor of a class is simply listed in parentheses immediately after the class name. There is no special keyword like `base`.

#### 5.3.1. Initializing and Coding Classes

By convention, the first argument of any Python class method (the reference to the current instance) is called `self`. This argument fills the role of `this` in C++ or `this` in Java.

#### 5.3.2. Knowing When to Use `self` and `__init__`

Methods are optional, but when you define one, you must remember to explicitly call the ancestor's `__init__` method (if it defines one). This is done by calling `super().__init__()`.

### 5.4. Instantiating Classes

In Python, simply call a class as if it were a function to create a new instance of the class. There is no explicit `new` operator like C++ or Java.

### 5.5. Exploring `UserDict`: A Wrapper Class

In the ActivePython IDE on Windows, you can quickly open any module in your library path by selecting File->Locate... (**Ctrl-L**).

Java and Powerbuilder support function overloading by argument list, *i.e.* one class can have multiple methods with the same name.

Guido, the original author of Python, explains method overriding this way: "Derived classes may override methods of their base class."

Always assign an initial value to all of an instance's data attributes in the `__init__` method. It will save you hours of debugging later, tracking down uninitialized attributes.

In versions of Python prior to 2.2, you could not directly subclass built-in datatypes like strings, lists, and dictionaries. To compensate, you had to use a wrapper class.

### 5.6.1. Getting and Setting Items

When accessing data attributes within a class, you need to qualify the attribute name: `obj.attr`. When calling other methods within a class, you use `obj.method()`.

### 5.7. Advanced Special Class Methods

In Java, you determine whether two string variables reference the same physical memory location by using `String.equals()`. This is called *object identity*.

While other object-oriented languages only let you define the physical model of an object ("this object has a `method`"), Python's `__str__` method lets you define the logical model.

### 5.8. Introducing Class Attributes

In Java, both static variables (called class attributes in Python) and instance variables (called data attributes in Python) are defined in the class.

There are no constants in Python. Everything can be changed if you try hard enough. This fits with one of the core principles of Python: *everything is mutable*.

### 5.9. Private Functions

In Python, all special methods (like `__str__`) and built-in attributes (like `__name__`) follow a standard naming convention: they both start with and end with two underscores.

## Chapter 6. Exceptions and File Handling

### 6.1. Handling Exceptions

Python uses `try/except` to handle exceptions and `raise` to generate them. Java and C++ use `try/catch` to handle exceptions, and `throw` to generate them.

### 6.5. Working with Directories

Whenever possible, you should use the functions in `os` and `shutil` for file, directory, and path manipulations. These modules are wrappers for the `os` and `shell` commands.

## Chapter 7. Regular Expressions

### 7.4. Using the `{n,m}` Syntax

There is no way to programmatically determine that two regular expressions are equivalent. The best you can do is write a lot of tests.

## Chapter 8. HTML Processing

### 8.2. Introducing `sgmlib.py`

Python 2.0 had a bug where `SGMLParser` would not recognize declarations at all (`!DOCTYPE` would never be called), which meant that `SGMLParser` were silently ignored.

In the ActivePython IDE on Windows, you can specify command line arguments in the "Run script" dialog. Separate multiple arguments with spaces.

### 8.4. Introducing `BaseHTMLProcessor.py`

The HTML specification requires that all non-HTML (like client-side JavaScript) must be enclosed in HTML comments, but not all HTML is valid.

### 8.5. locals and globals

Python 2.2 introduced a subtle but important change that affects the namespace search order: nested scopes. In versions of Python prior to 2.2, the search order was different.

Using the `getattr` and `setattr` functions, you can get the value of arbitrary variables dynamically, providing the variable name as a string. This mirrors the behavior of the `getattr` and `setattr` functions in the `os` module.

Using dictionary-based string formatting with `%` is a convenient way of making complex string formatting expressions more readable.

## Chapter 9. XML Processing

### 9.2. Packages

A package is a directory with the special `__init__.py` file in it. The `__init__.py` file defines the attributes and methods of the package. It doesn't need to define anything.

### 9.6. Accessing element attributes

This section may be a little confusing, because of some overlapping terminology. Elements in an XML document have attributes, and attributes have values.

Like a dictionary, attributes of an XML element have no ordering. Attributes may *happen to be* listed in a certain order in the original document, but this is not guaranteed.

## Chapter 10. Scripts and Streams

## Chapter 11. HTTP Web Services

### 11.6. Handling Last-Modified and ETag

In these examples, the HTTP server has supported both `ETag` and `Last-Modified` headers, but not all servers do. As a web services client, you should be able to handle either.

## Chapter 12. SOAP Web Services

## Chapter 13. Unit Testing

### 13.2. Diving in

`unittest` is included with Python 2.1 and later. Python 2.0 users can download it from <http://pyunit.sourceforge.net/>.

## Chapter 14. Test-First Programming

### 14.3. `roman.py`, stage 3

The most important thing that comprehensive unit testing can tell you is when to stop coding. When all the unit tests for a function pass, you can stop coding.

### 14.5. `roman.py`, stage 5

When all of your tests pass, stop coding.

## Chapter 15. Refactoring

### 15.3. Refactoring

Whenever you are going to use a regular expression more than once, you should compile it to get a pattern object, then call the method

## Chapter 16. Functional Programming

### 16.2. Finding the path

The pathnames and filenames you pass to `os.path` do not need to exist.

`os.path` not only constructs full path names, it also normalizes them. That means that if you are in the `current` directory, `os.path` will return `os.path`. It normalizes

Like the other functions in the `os` and `os.path` modules, `os.path` is cross-platform. Your results will look slightly different than my examples if you're

## Chapter 17. Dynamic functions

## Chapter 18. Performance Tuning

### 18.2. Using the `timeit` Module

You can use the `timeit` module on the command line to test an existing Python program, without modifying the code. See <http://docs.python.org/2/library/timeit.html>

The `timeit` module only works if you already know what piece of code you need to optimize. If you have a larger Python program and don't

# Appendix D. List of examples

## Chapter 1. Installing Python

### 1.3. Python on Mac OS X

Example 1.1. Two versions of Python

### 1.5. Python on RedHat Linux

Example 1.2. Installing on RedHat Linux 9

### 1.6. Python on Debian GNU/Linux

Example 1.3. Installing on Debian GNU/Linux

### 1.7. Python Installation from Source

Example 1.4. Installing from source

### 1.8. The Interactive Shell

Example 1.5. First Steps in the Interactive Shell

## Chapter 2. Your First Python Program

### 2.1. Diving in

Example 2.1. odbchelper.py

### 2.3. Documenting Functions

Example 2.2. Defining the buildConnectionString Function's doc string

### 2.4. Everything Is an Object

Example 2.3. Accessing the buildConnectionString Function's doc string

#### 2.4.1. The Import Search Path

Example 2.4. Import Search Path

### 2.5. Indenting Code

Example 2.5. Indenting the buildConnectionString Function

Example 2.6. if Statements

## Chapter 3. Native Datatypes

### 3.1.1. Defining Dictionaries

Example 3.1. Defining a Dictionary

### 3.1.2. Modifying Dictionaries

Example 3.2. Modifying a Dictionary

Example 3.3. Dictionary Keys Are Case-Sensitive

Example 3.4. Mixing Datatypes in a Dictionary

### 3.1.3. Deleting Items From Dictionaries

Example 3.5. Deleting Items from a Dictionary

### 3.2.1. Defining Lists

Example 3.6. Defining a List

Example 3.7. Negative List Indices

Example 3.8. Slicing a List

Example 3.9. Slicing Shorthand

### 3.2.2. Adding Elements to Lists

Example 3.10. Adding Elements to a List

Example 3.11. The Difference between extend and append

### 3.2.3. Searching Lists

Example 3.12. Searching a List

### 3.2.4. Deleting List Elements

Example 3.13. Removing Elements from a List

### 3.2.5. Using List Operators

Example 3.14. List Operators

- 3.3. Introducing Tuples
  - Example 3.15. Defining a tuple
  - Example 3.16. Tuples Have No Methods
- 3.4. Declaring variables
  - Example 3.17. Defining the myParams Variable
- 3.4.1. Referencing Variables
  - Example 3.18. Referencing an Unbound Variable
- 3.4.2. Assigning Multiple Values at Once
  - Example 3.19. Assigning multiple values at once
  - Example 3.20. Assigning Consecutive Values
- 3.5. Formatting Strings
  - Example 3.21. Introducing String Formatting
  - Example 3.22. String Formatting vs. Concatenating
  - Example 3.23. Formatting Numbers
- 3.6. Mapping Lists
  - Example 3.24. Introducing List Comprehensions
  - Example 3.25. The keys, values, and items Functions
  - Example 3.26. List Comprehensions in buildConnectionString, Step by Step
- 3.7. Joining Lists and Splitting Strings
  - Example 3.27. Output of odbchelper.py
  - Example 3.28. Splitting a String

## Chapter 4. The Power Of Introspection

- 4.1. Diving In
  - Example 4.1. apihelper.py
  - Example 4.2. Sample Usage of apihelper.py
  - Example 4.3. Advanced Usage of apihelper.py
- 4.2. Using Optional and Named Arguments
  - Example 4.4. Valid Calls of info
- 4.3.1. The type Function
  - Example 4.5. Introducing type
- 4.3.2. The str Function
  - Example 4.6. Introducing str
  - Example 4.7. Introducing dir
  - Example 4.8. Introducing callable
- 4.3.3. Built-In Functions
  - Example 4.9. Built-in Attributes and Functions
- 4.4. Getting Object References With getattr
  - Example 4.10. Introducing getattr
- 4.4.1. getattr with Modules
  - Example 4.11. The getattr Function in apihelper.py
- 4.4.2. getattr As a Dispatcher
  - Example 4.12. Creating a Dispatcher with getattr
  - Example 4.13. getattr Default Values
- 4.5. Filtering Lists
  - Example 4.14. Introducing List Filtering
- 4.6. The Peculiar Nature of and and or
  - Example 4.15. Introducing and
  - Example 4.16. Introducing or
- 4.6.1. Using the and-or Trick
  - Example 4.17. Introducing the and-or Trick
  - Example 4.18. When the and-or Trick Fails

- Example 4.19. Using the and-or Trick Safely
- 4.7. Using lambda Functions
  - Example 4.20. Introducing lambda Functions
- 4.7.1. Real-World lambda Functions
  - Example 4.21. split With No Arguments
- 4.8. Putting It All Together
  - Example 4.22. Getting a doc string Dynamically
  - Example 4.23. Why Use str on a doc string?
  - Example 4.24. Introducing ljust
  - Example 4.25. Printing a List

## Chapter 5. Objects and Object-Orientation

- 5.1. Diving In
  - Example 5.1. fileinfo.py
- 5.2. Importing Modules Using from module import
  - Example 5.2. import module vs. from module import
- 5.3. Defining Classes
  - Example 5.3. The Simplest Python Class
  - Example 5.4. Defining the FileInfo Class
- 5.3.1. Initializing and Coding Classes
  - Example 5.5. Initializing the FileInfo Class
  - Example 5.6. Coding the FileInfo Class
- 5.4. Instantiating Classes
  - Example 5.7. Creating a FileInfo Instance
- 5.4.1. Garbage Collection
  - Example 5.8. Trying to Implement a Memory Leak
- 5.5. Exploring UserDict: A Wrapper Class
  - Example 5.9. Defining the UserDict Class
  - Example 5.10. UserDict Normal Methods
  - Example 5.11. Inheriting Directly from Built-In Datatype dict
- 5.6.1. Getting and Setting Items
  - Example 5.12. The \_\_getitem\_\_ Special Method
  - Example 5.13. The \_\_setitem\_\_ Special Method
  - Example 5.14. Overriding \_\_setitem\_\_ in MP3FileInfo
  - Example 5.15. Setting an MP3FileInfo's name
- 5.7. Advanced Special Class Methods
  - Example 5.16. More Special Methods in UserDict
- 5.8. Introducing Class Attributes
  - Example 5.17. Introducing Class Attributes
  - Example 5.18. Modifying Class Attributes
- 5.9. Private Functions
  - Example 5.19. Trying to Call a Private Method

## Chapter 6. Exceptions and File Handling

- 6.1. Handling Exceptions
  - Example 6.1. Opening a Non-Existent File
- 6.1.1. Using Exceptions For Other Purposes
  - Example 6.2. Supporting Platform-Specific Functionality
- 6.2. Working with File Objects
  - Example 6.3. Opening a File
- 6.2.1. Reading Files



- Example 6.4. Reading a File
- 6.2.2. Closing Files
  - Example 6.5. Closing a File
- 6.2.3. Handling I/O Errors
  - Example 6.6. File Objects in MP3FileInfo
- 6.2.4. Writing to Files
  - Example 6.7. Writing to Files
- 6.3. Iterating with for Loops
  - Example 6.8. Introducing the for Loop
  - Example 6.9. Simple Counters
  - Example 6.10. Iterating Through a Dictionary
  - Example 6.11. for Loop in MP3FileInfo
- 6.4. Using sys.modules
  - Example 6.12. Introducing sys.modules
  - Example 6.13. Using sys.modules
  - Example 6.14. The \_\_module\_\_ Class Attribute
  - Example 6.15. sys.modules in fileinfo.py
- 6.5. Working with Directories
  - Example 6.16. Constructing Pathnames
  - Example 6.17. Splitting Pathnames
  - Example 6.18. Listing Directories
  - Example 6.19. Listing Directories in fileinfo.py
  - Example 6.20. Listing Directories with glob
- 6.6. Putting It All Together
  - Example 6.21. listDirectory

## Chapter 7. Regular Expressions

- 7.2. Case Study: Street Addresses
  - Example 7.1. Matching at the End of a String
  - Example 7.2. Matching Whole Words
- 7.3.1. Checking for Thousands
  - Example 7.3. Checking for Thousands
- 7.3.2. Checking for Hundreds
  - Example 7.4. Checking for Hundreds
- 7.4. Using the {n,m} Syntax
  - Example 7.5. The Old Way: Every Character Optional
  - Example 7.6. The New Way: From n o m
- 7.4.1. Checking for Tens and Ones
  - Example 7.7. Checking for Tens
  - Example 7.8. Validating Roman Numerals with {n,m}
- 7.5. Verbose Regular Expressions
  - Example 7.9. Regular Expressions with Inline Comments
- 7.6. Case study: Parsing Phone Numbers
  - Example 7.10. Finding Numbers
  - Example 7.11. Finding the Extension
  - Example 7.12. Handling Different Separators
  - Example 7.13. Handling Numbers Without Separators
  - Example 7.14. Handling Leading Characters
  - Example 7.15. Phone Number, Wherever I May Find Ye
  - Example 7.16. Parsing Phone Numbers (Final Version)

## Chapter 8. HTML Processing

- 8.1. Diving in
  - Example 8.1. BaseHTMLProcessor.py
  - Example 8.2. dialect.py
  - Example 8.3. Output of dialect.py
- 8.2. Introducing sgmlib.py
  - Example 8.4. Sample test of sgmlib.py
- 8.3. Extracting data from HTML documents
  - Example 8.5. Introducing urllib
  - Example 8.6. Introducing urlister.py
  - Example 8.7. Using urlister.py
- 8.4. Introducing BaseHTMLProcessor.py
  - Example 8.8. Introducing BaseHTMLProcessor
  - Example 8.9. BaseHTMLProcessor output
- 8.5. locals and globals
  - Example 8.10. Introducing locals
  - Example 8.11. Introducing globals
  - Example 8.12. locals is read-only, globals is not
- 8.6. Dictionary-based string formatting
  - Example 8.13. Introducing dictionary-based string formatting
  - Example 8.14. Dictionary-based string formatting in BaseHTMLProcessor.py
  - Example 8.15. More dictionary-based string formatting
- 8.7. Quoting attribute values
  - Example 8.16. Quoting attribute values
- 8.8. Introducing dialect.py
  - Example 8.17. Handling specific tags
  - Example 8.18. SGMLParser
  - Example 8.19. Overriding the handle\_data method
- 8.9. Putting it all together
  - Example 8.20. The translate function, part 1
  - Example 8.21. The translate function, part 2: curiouser and curiouser
  - Example 8.22. The translate function, part 3

## Chapter 9. XML Processing

- 9.1. Diving in
  - Example 9.1. kgp.py
  - Example 9.2. toolbox.py
  - Example 9.3. Sample output of kgp.py
  - Example 9.4. Simpler output from kgp.py
- 9.2. Packages
  - Example 9.5. Loading an XML document (a sneak peek)
  - Example 9.6. File layout of a package
  - Example 9.7. Packages are modules, too
- 9.3. Parsing XML
  - Example 9.8. Loading an XML document (for real this time)
  - Example 9.9. Getting child nodes
  - Example 9.10. toxml works on any node
  - Example 9.11. Child nodes can be text
  - Example 9.12. Drilling down all the way to text
- 9.4. Unicode
  - Example 9.13. Introducing unicode
  - Example 9.14. Storing non-ASCII characters
  - Example 9.15. sitecustomize.py

Example 9.16. Effects of setting the default encoding

Example 9.17. Specifying encoding in .py files

Example 9.18. russiansample.xml

Example 9.19. Parsing russiansample.xml

#### 9.5. Searching for elements

Example 9.20. binary.xml

Example 9.21. Introducing getElementByTagName

Example 9.22. Every element is searchable

Example 9.23. Searching is actually recursive

#### 9.6. Accessing element attributes

Example 9.24. Accessing element attributes

Example 9.25. Accessing individual attributes

### Chapter 10. Scripts and Streams

#### 10.1. Abstracting input sources

Example 10.1. Parsing XML from a file

Example 10.2. Parsing XML from a URL

Example 10.3. Parsing XML from a string (the easy but inflexible way)

Example 10.4. Introducing StringIO

Example 10.5. Parsing XML from a string (the file-like object way)

Example 10.6. openAnything

Example 10.7. Using openAnything in kgp.py

#### 10.2. Standard input, output, and error

Example 10.8. Introducing stdout and stderr

Example 10.9. Redirecting output

Example 10.10. Redirecting error information

Example 10.11. Printing to stderr

Example 10.12. Chaining commands

Example 10.13. Reading from standard input in kgp.py

#### 10.3. Caching node lookups

Example 10.14. loadGrammar

Example 10.15. Using the ref element cache

#### 10.4. Finding direct children of a node

Example 10.16. Finding direct child elements

#### 10.5. Creating separate handlers by node type

Example 10.17. Class names of parsed XML objects

Example 10.18. parse, a generic XML node dispatcher

Example 10.19. Functions called by the parse dispatcher

#### 10.6. Handling command-line arguments

Example 10.20. Introducing sys.argv

Example 10.21. The contents of sys.argv

Example 10.22. Introducing getopt

Example 10.23. Handling command-line arguments in kgp.py

### Chapter 11. HTTP Web Services

#### 11.1. Diving in

Example 11.1. openanything.py

#### 11.2. How not to fetch data over HTTP

Example 11.2. Downloading a feed the quick-and-dirty way

#### 11.4. Debugging HTTP web services

Example 11.3. Debugging HTTP

## 11.5. Setting the User-Agent

Example 11.4. Introducing urllib2

Example 11.5. Adding headers with the Request

## 11.6. Handling Last-Modified and ETag

Example 11.6. Testing Last-Modified

Example 11.7. Defining URL handlers

Example 11.8. Using custom URL handlers

Example 11.9. Supporting ETag/If-None-Match

## 11.7. Handling redirects

Example 11.10. Accessing web services without a redirect handler

Example 11.11. Defining the redirect handler

Example 11.12. Using the redirect handler to detect permanent redirects

Example 11.13. Using the redirect handler to detect temporary redirects

## 11.8. Handling compressed data

Example 11.14. Telling the server you would like compressed data

Example 11.15. Decompressing the data

Example 11.16. Decompressing the data directly from the server

## 11.9. Putting it all together

Example 11.17. The openanything function

Example 11.18. The fetch function

Example 11.19. Using openanything.py

# Chapter 12. SOAP Web Services

## 12.1. Diving In

Example 12.1. search.py

Example 12.2. Sample Usage of search.py

## 12.2.1. Installing PyXML

Example 12.3. Verifying PyXML Installation

## 12.2.2. Installing fpconst

Example 12.4. Verifying fpconst Installation

## 12.2.3. Installing SOAPpy

Example 12.5. Verifying SOAPpy Installation

## 12.3. First Steps with SOAP

Example 12.6. Getting the Current Temperature

## 12.4. Debugging SOAP Web Services

Example 12.7. Debugging SOAP Web Services

## 12.6. Introspecting SOAP Web Services with WSDL

Example 12.8. Discovering The Available Methods

Example 12.9. Discovering A Method's Arguments

Example 12.10. Discovering A Method's Return Values

Example 12.11. Calling A Web Service Through A WSDL Proxy

## 12.7. Searching Google

Example 12.12. Introspecting Google Web Services

Example 12.13. Searching Google

Example 12.14. Accessing Secondary Information From Google

## 12.8. Troubleshooting SOAP Web Services

Example 12.15. Calling a Method With an Incorrectly Configured Proxy

Example 12.16. Calling a Method With the Wrong Arguments

Example 12.17. Calling a Method and Expecting the Wrong Number of Return Values

Example 12.18. Calling a Method With An Application-Specific Error

# Chapter 13. Unit Testing

- 13.3. Introducing `romantest.py`
  - Example 13.1. `romantest.py`
- 13.4. Testing for success
  - Example 13.2. `testToRomanKnownValues`
- 13.5. Testing for failure
  - Example 13.3. Testing bad input to `toRoman`
  - Example 13.4. Testing bad input to `fromRoman`
- 13.6. Testing for sanity
  - Example 13.5. Testing `toRoman` against `fromRoman`
  - Example 13.6. Testing for case

## Chapter 14. Test-First Programming

- 14.1. `roman.py`, stage 1
  - Example 14.1. `roman1.py`
  - Example 14.2. Output of `romantest1.py` against `roman1.py`
- 14.2. `roman.py`, stage 2
  - Example 14.3. `roman2.py`
  - Example 14.4. How `toRoman` works
  - Example 14.5. Output of `romantest2.py` against `roman2.py`
- 14.3. `roman.py`, stage 3
  - Example 14.6. `roman3.py`
  - Example 14.7. Watching `toRoman` handle bad input
  - Example 14.8. Output of `romantest3.py` against `roman3.py`
- 14.4. `roman.py`, stage 4
  - Example 14.9. `roman4.py`
  - Example 14.10. How `fromRoman` works
  - Example 14.11. Output of `romantest4.py` against `roman4.py`
- 14.5. `roman.py`, stage 5
  - Example 14.12. `roman5.py`
  - Example 14.13. Output of `romantest5.py` against `roman5.py`

## Chapter 15. Refactoring

- 15.1. Handling bugs
  - Example 15.1. The bug
  - Example 15.2. Testing for the bug (`romantest61.py`)
  - Example 15.3. Output of `romantest61.py` against `roman61.py`
  - Example 15.4. Fixing the bug (`roman62.py`)
  - Example 15.5. Output of `romantest62.py` against `roman62.py`
- 15.2. Handling changing requirements
  - Example 15.6. Modifying test cases for new requirements (`romantest71.py`)
  - Example 15.7. Output of `romantest71.py` against `roman71.py`
  - Example 15.8. Coding the new requirements (`roman72.py`)
  - Example 15.9. Output of `romantest72.py` against `roman72.py`
- 15.3. Refactoring
  - Example 15.10. Compiling regular expressions
  - Example 15.11. Compiled regular expressions in `roman81.py`
  - Example 15.12. Output of `romantest81.py` against `roman81.py`
  - Example 15.13. `roman82.py`
  - Example 15.14. Output of `romantest82.py` against `roman82.py`
  - Example 15.15. `roman83.py`
  - Example 15.16. Output of `romantest83.py` against `roman83.py`

#### 15.4. Postscript

Example 15.17. `roman9.py`

Example 15.18. Output of `romantest9.py` against `roman9.py`

### Chapter 16. Functional Programming

#### 16.1. Diving in

Example 16.1. `regression.py`

Example 16.2. Sample output of `regression.py`

#### 16.2. Finding the path

Example 16.3. `fullpath.py`

Example 16.4. Further explanation of `os.path.abspath`

Example 16.5. Sample output from `fullpath.py`

Example 16.6. Running scripts in the current directory

#### 16.3. Filtering lists revisited

Example 16.7. Introducing filter

Example 16.8. filter in `regression.py`

Example 16.9. Filtering using list comprehensions instead

#### 16.4. Mapping lists revisited

Example 16.10. Introducing map

Example 16.11. map with lists of mixed datatypes

Example 16.12. map in `regression.py`

#### 16.6. Dynamically importing modules

Example 16.13. Importing multiple modules at once

Example 16.14. Importing modules dynamically

Example 16.15. Importing a list of modules dynamically

#### 16.7. Putting it all together

Example 16.16. The `regressionTest` function

Example 16.17. Step 1: Get all the files

Example 16.18. Step 2: Filter to find the files you care about

Example 16.19. Step 3: Map filenames to module names

Example 16.20. Step 4: Mapping module names to modules

Example 16.21. Step 5: Loading the modules into a test suite

Example 16.22. Step 6: Telling unittest to use your test suite

### Chapter 17. Dynamic functions

#### 17.2. `plural.py`, stage 1

Example 17.1. `plural1.py`

Example 17.2. Introducing `re.sub`

Example 17.3. Back to `plural1.py`

Example 17.4. More on negation regular expressions

Example 17.5. More on `re.sub`

#### 17.3. `plural.py`, stage 2

Example 17.6. `plural2.py`

Example 17.7. Unrolling the plural function

#### 17.4. `plural.py`, stage 3

Example 17.8. `plural3.py`

#### 17.5. `plural.py`, stage 4

Example 17.9. `plural4.py`

Example 17.10. `plural4.py` continued

Example 17.11. Unrolling the rules definition

Example 17.12. `plural4.py`, finishing up

- Example 17.13. Another look at buildMatchAndApplyFunctions
- Example 17.14. Expanding tuples when calling functions
- 17.6. plural.py, stage 5
  - Example 17.15. rules.en
  - Example 17.16. plural5.py
- 17.7. plural.py, stage 6
  - Example 17.17. plural6.py
  - Example 17.18. Introducing generators
  - Example 17.19. Using generators instead of recursion
  - Example 17.20. Generators in for loops
  - Example 17.21. Generators that generate dynamic functions

## Chapter 18. Performance Tuning

- 18.1. Diving in
  - Example 18.1. soundex/stage1/soundex1a.py
- 18.2. Using the timeit Module
  - Example 18.2. Introducing timeit
- 18.3. Optimizing Regular Expressions
  - Example 18.3. Best Result So Far: soundex/stage1/soundex1e.py
- 18.4. Optimizing Dictionary Lookups
  - Example 18.4. Best Result So Far: soundex/stage2/soundex2c.py
- 18.5. Optimizing List Operations
  - Example 18.5. Best Result So Far: soundex/stage2/soundex2c.py

## Appendix E. Revision history

Revision History	
Revision 5.4	2004-05-20
Added Section 12.1, “Diving In”. Added Section 12.2, “Installing the SOAP Libraries”. Added Section 12.3, “First Steps with SOAP”. Added Section 12.4, “Debugging SOAP Web Services”. Added Section 12.5, “Introducing WSDL”. Added Section 12.6, “Introspecting SOAP Web Services with WSDL”. Added Section 12.7, “Searching Google”. Added Section 12.8, “Troubleshooting SOAP Web Services”. Added Section 12.9, “Summary”. Incorporated technical reviewer revisions in Chapter 16, <i>Functional Programming</i> and Chapter 18, <i>Performance Tuning</i> .	
Revision 5.3	2004-05-12
Added Example to Section 18.3, “Optimizing Regular Expressions”. Thanks, Paul. Incorporated copyediting revisions into Chapter 5, <i>Objects and Object-Oriented</i> and Chapter 6, <i>Exceptions and File Handling</i> . Fixed URL of Section 9.7, “Segue”.	
Revision 5.2	2004-05-09
Fixed URL of Section 14.1, “roman.py, stage 1”. Added Section 18.1, “Diving in”. Added Section 18.2, “Using the timeit Module”. Added Section 18.3, “Optimizing Regular Expressions”. Added Section 18.4, “Optimizing Dictionary Lookups”. Added Section 18.5, “Optimizing List Operations”. Added Section 18.6, “Optimizing String Manipulation”. Added Section 18.7, “Summary”.	
Revision 5.1	2004-05-05
Clarified Example 7.7, “Checking for Tens” and Example 7.8, “Validating Roman Numerals with {n,m}”. Clarified Example 7.10, “Finding Numbers”. Fixed typo in Example 11.6, “Testing Last-Modified”. Thanks, Jesir. Fixed typo in Example 3.11, “The Difference between extend and append”. Thanks, Daniel. Incorporated technical reviewer revisions. Incorporated copy editor revisions in Chapter 1, <i>Installing Python</i> , Chapter 2, <i>Your First Python Program</i> , Chapter 3, <i>Native Data</i>	
Revision 5.0	2004-04-16
Added Section 11.1, “Diving in”. Added Section 11.2, “How not to fetch data over HTTP”. Added Section 11.3, “Features of HTTP”. Added Section 11.4, “Debugging HTTP web services”. Added Section 11.5, “Setting the User-Agent”. Added Section 11.6, “Handling Last-Modified and ETag”. Added Section 11.7, “Handling redirects”. Added Section 11.8, “Handling compressed data”. Added Section 11.9, “Putting it all together”. Added Section 11.10, “Summary”. Added Example 3.11, “The Difference between extend and append”. Changed descriptions of how to download Python throughout Chapter 1, <i>Installing Python</i> to be more generic and less version-specific. Changed references of “module” to “program” in Section 2.1, “Diving in” and Section 2.4, “Everything Is an Object” since we have Added explicit instructions in Section 2.4, “Everything Is an Object” for the reader to open their Python IDE and follow along with	



<p>Changed all examples and descriptions that referred to truth values <code>1</code> and <code>0</code> to refer to <code>True</code> and <code>False</code>.</p> <p>Updated Example 3.22, “String Formatting vs. Concatenating” to show new Python 2.3 <code>message</code>.</p> <p>Fixed typo in Example 17.19, “Using generators instead of recursion”.</p> <p>Fixed typo in Section 7.7, “Summary”.</p> <p>Fixed typo in Example 17.9, “plural4.py”.</p>	
Revision 4.9	2004-03-25
<p>Finished Section 16.7, “Putting it all together”.</p> <p>Added Section 16.8, “Summary”.</p> <p>Split unit testing introduction into two chapters, Chapter 13, <i>Unit Testing</i> and Chapter 14, <i>Test-First Programming</i>.</p> <p>Fixed typo in Example 17.12, “plural4.py, finishing up”.</p> <p>Fixed typo in Example 17.18, “Introducing generators”.</p>	
Revision 4.8	2004-03-25
<p>Finished Section 17.7, “plural.py, stage 6”.</p> <p>Finished Section 17.8, “Summary”.</p> <p>Fixed broken links in Appendix A, <i>Further reading</i>, Appendix B, <i>A 5-minute review</i>, Appendix C, <i>Tips and tricks</i>, Appendix D, <i>List of...</i></p>	
Revision 4.7	2004-03-21
<p>Added Section 17.1, “Diving in”.</p> <p>Added Section 17.2, “plural.py, stage 1”.</p> <p>Added Section 17.3, “plural.py, stage 2”.</p> <p>Added Section 17.4, “plural.py, stage 3”.</p> <p>Added Section 17.5, “plural.py, stage 4”.</p> <p>Added Section 17.6, “plural.py, stage 5”.</p> <p>Added Section 17.7, “plural.py, stage 6” (unfinished).</p> <p>Added Section 17.8, “Summary” (unfinished).</p>	
Revision 4.6	2004-03-14
<p>Finished Section 7.4, “Using the {n,m} Syntax”.</p> <p>Finished Section 7.5, “Verbose Regular Expressions”.</p> <p>Finished Section 7.6, “Case study: Parsing Phone Numbers”.</p> <p>Expanded Section 7.7, “Summary”.</p>	
Revision 4.5	2004-03-07
<p>Added Section 7.1, “Diving In”.</p> <p>Added Section 7.4, “Using the {n,m} Syntax” (incomplete).</p> <p>Added Section 7.5, “Verbose Regular Expressions” (incomplete).</p> <p>Added Section 7.6, “Case study: Parsing Phone Numbers” (incomplete).</p> <p>Added Section 7.7, “Summary”.</p> <p>Moved Section 7.2, “Case Study: Street Addresses” and Section 7.3, “Case Study: Roman Numerals” to regular expressions chapter.</p> <p>Added Example 6.20, “Listing Directories with glob”.</p> <p>Added Example 6.7, “Writing to Files”.</p> <p>Added Example 5.11, “Inheriting Directly from Built-In Datatype dict”.</p> <p>Added Example 10.11, “Printing to stderr”.</p> <p>Added Example 4.12, “Creating a Dispatcher with getattr” and Example 4.13, “getattr Default Values”.</p> <p>Added Example 2.6, “if Statements”.</p> <p>Added Example 3.23, “Formatting Numbers”.</p> <p>Split Chapter 5, <i>Objects and Object-Orientation</i> into 2 chapters: Chapter 5, <i>Objects and Object-Orientation</i> and Chapter 6, <i>Exceptions and Error Handling</i>.</p> <p>Split Chapter 9, <i>XML Processing</i> into 2 chapters: Chapter 9, <i>XML Processing</i> and Chapter 10, <i>Scripts and Streams</i>.</p> <p>Split Chapter 13, <i>Unit Testing</i> into 2 chapters: Chapter 13, <i>Unit Testing</i> and Chapter 15, <i>Refactoring</i>.</p> <p>Renamed <code>__doc__</code> in Chapter 4, <i>The Power Of Introspection</i>.</p> <p>Fixed incorrect back-reference in Section 8.5, “locals and globals”.</p> <p>Fixed broken example links in Section 8.1, “Diving in”.</p> <p>Fixed missing line in example in Section 9.1, “Diving in”.</p>	

Fixed typo in Section 8.2, “Introducing sgmlib.py”.	
Revision 4.4	2003-10-08
<p>Added Section 1.1, “Which Python is right for you?”.</p> <p>Added Section 1.2, “Python on Windows”.</p> <p>Added Section 1.3, “Python on Mac OS X”.</p> <p>Added Section 1.4, “Python on Mac OS 9”.</p> <p>Added Section 1.5, “Python on RedHat Linux”.</p> <p>Added Section 1.6, “Python on Debian GNU/Linux”.</p> <p>Added Section 1.7, “Python Installation from Source”.</p> <p>Added Section 1.9, “Summary”.</p> <p>Removed preface.</p> <p>Fixed typo in Example 3.27, “Output of odbchelper.py”.</p> <p>Added link to PEP 257 in Section 2.3, “Documenting Functions”.</p> <p>Fixed link to <i>How to Think Like a Computer Scientist</i> (<a href="http://www.ibiblio.org/obp/thinkCSpy/">http://www.ibiblio.org/obp/thinkCSpy/</a>) in Section 3.4.2, “Assigning Multiple Values”.</p> <p>Added note about implied assert in Section 3.3, “Introducing Tuples”.</p>	
Revision 4.3	2003-09-28
<p>Added Section 16.6, “Dynamically importing modules”.</p> <p>Added Section 16.7, “Putting it all together” (incomplete).</p> <p>Fixed links in Appendix F, <i>About the book</i>.</p>	
Revision 4.2.1	2003-09-17
<p>Fixed links on index page.</p> <p>Fixed syntax highlighting.</p>	
Revision 4.2	2003-09-12
<p>Fixed typos in Section 16.4, “Mapping lists revisited”, Section 16.3, “Filtering lists revisited”, Section 7.2, “Case Study: Street Address Parsing”.</p> <p>Fixed external link in Section 5.3, “Defining Classes”. Thanks, Harry.</p> <p>Changed wording at the end of Section 4.5, “Filtering Lists”. Thanks, Paul.</p> <p>Added sentence in Section 13.5, “Testing for failure” to make it clearer that we’re passing a function to <code>assert</code>, not a function name as a string.</p> <p>Fixed typo in Section 8.8, “Introducing dialect.py”. Thanks, Wellie.</p> <p>Fixed links to dialectized examples.</p> <p>Fixed external link to the history of Roman numerals. Thanks to many concerned Roman numeral fans around the world.</p>	
Revision 4.1	2002-07-28
<p>Added Section 10.3, “Caching node lookups”.</p> <p>Added Section 10.4, “Finding direct children of a node”.</p> <p>Added Section 10.5, “Creating separate handlers by node type”.</p> <p>Added Section 10.6, “Handling command-line arguments”.</p> <p>Added Section 10.7, “Putting it all together”.</p> <p>Added Section 10.8, “Summary”.</p> <p>Fixed typo in Section 6.5, “Working with Directories”. It’s <code>os.listdir()</code>, not <code>os.listdirs()</code>. Thanks, Abhishek.</p> <p>Fixed typo in Section 3.7, “Joining Lists and Splitting Strings”. When evaluated (instead of printed), the Python IDE will display single quotes around strings.</p> <p>Changed example in Section 4.8, “Putting It All Together” to use a user-defined function, since Python 2.2 obsoleted the old example.</p> <p>Fixed typo in Section 9.4, “Unicode”, “anyway” to “anywhere”. Thanks Frank.</p> <p>Fixed typo in Section 13.6, “Testing for sanity”, doubled word “accept”. Thanks Ralph.</p> <p>Fixed typo in Section 15.3, “Refactoring”, “matches 0 to 3 characters, not 4”. Thanks Ralph.</p> <p>Clarified and expanded explanation of implied slice indices in Example 3.9, “Slicing Shorthand”. Thanks Petr.</p> <p>Added historical note in Section 5.5, “Exploring UserDict: A Wrapper Class” now that Python 2.2 supports subclassing built-in data types.</p> <p>Added explanation of <code>dict.items()</code> method in Example 5.9, “Defining the UserDict Class”. Thanks Petr.</p> <p>Clarified Python’s lack of overloading in Section 5.5, “Exploring UserDict: A Wrapper Class”. Thanks Petr.</p> <p>Fixed typo in Example 8.8, “Introducing BaseHTMLProcessor”. HTML comments end with two dashes and a bracket, not one. Thanks Petr.</p> <p>Changed tense of note about nested scopes in Section 8.5, “locals and globals” now that Python 2.2 is out. Thanks Petr.</p> <p>Fixed typo in Example 8.14, “Dictionary-based string formatting in BaseHTMLProcessor.py”; a space should have been a non-breaking space.</p>	

<p>Added title to note on derived classes in Section 5.5, “Exploring UserDict: A Wrapper Class”. Thanks Petr.</p> <p>Added title to note on downloading <code>urllib</code> in Section 15.3, “Refactoring”. Thanks Petr.</p> <p>Fixed typesetting problem in Example 16.6, “Running scripts in the current directory”; tabs should have been spaces, and the line number should have been 16.</p> <p>Fixed capitalization typo in the tip on truth values in Section 3.2, “Introducing Lists”. It’s <code>and</code> not <code>and</code>. Thanks to everyone who pointed it out.</p> <p>Changed section titles of Section 3.1, “Introducing Dictionaries”, Section 3.2, “Introducing Lists”, and Section 3.3, “Introducing Tuples”.</p> <p>Upgraded to version 1.52 of the DocBook XSL stylesheets.</p> <p>Upgraded to version 6.52 of the SAXON XSLT processor from Michael Kay.</p> <p>Various accessibility-related stylesheet tweaks.</p> <p>Somewhere between this revision and the last one, she said yes. The wedding will be next spring.</p>	
Revision 4.0-2	2002-04-26
<p>Fixed typo in Example 4.15, “Introducing and”.</p> <p>Fixed typo in Example 2.4, “Import Search Path”.</p> <p>Fixed Windows help file (missing table of contents due to base stylesheet changes).</p>	
Revision 4.0	2002-04-19
<p>Expanded Section 2.4, “Everything Is an Object” to include more about import search paths.</p> <p>Fixed typo in Example 3.7, “Negative List Indices”. Thanks to Brian for the correction.</p> <p>Rewrote the tip on truth values in Section 3.2, “Introducing Lists”, now that Python has a separate boolean datatype.</p> <p>Fixed typo in Section 5.2, “Importing Modules Using from module import” when comparing syntax to Java. Thanks to Rick for the correction.</p> <p>Added note in Section 5.5, “Exploring UserDict: A Wrapper Class” about derived classes always overriding ancestor classes.</p> <p>Fixed typo in Example 5.18, “Modifying Class Attributes”. Thanks to Kevin for the correction.</p> <p>Added note in Section 6.1, “Handling Exceptions” that you can define and raise your own exceptions. Thanks to Rony for the suggestion.</p> <p>Fixed typo in Example 8.17, “Handling specific tags”. Thanks for Rick for the correction.</p> <p>Added note in Example 8.18, “SGMLParser” about what the return codes mean. Thanks to Howard for the suggestion.</p> <p>Added <code>urllib</code> function when creating <code>urllib</code> instance in Example 10.6, “openAnything”. Thanks to Ganesan for the idea.</p> <p>Added link in Section 13.3, “Introducing romantest.py” to explanation of why test cases belong in a separate file.</p> <p>Changed Section 16.2, “Finding the path” to use <code>os.path</code> instead of <code>os</code>. Thanks to Marc for the idea.</p> <p>Added code samples (<code>re</code> and <code>re</code>) for the upcoming regular expressions chapter.</p> <p>Updated and expanded list of Python distributions on home page.</p>	
Revision 3.9	2002-01-01
<p>Added Section 9.4, “Unicode”.</p> <p>Added Section 9.5, “Searching for elements”.</p> <p>Added Section 9.6, “Accessing element attributes”.</p> <p>Added Section 10.1, “Abstracting input sources”.</p> <p>Added Section 10.2, “Standard input, output, and error”.</p> <p>Added simple counter <code>for</code> loop examples (good usage and bad usage) in Section 6.3, “Iterating with for Loops”. Thanks to Kevin for the suggestion.</p> <p>Fixed typo in Example 3.25, “The keys, values, and items Functions” (two elements of <code>items</code> were reversed).</p> <p>Fixed mistake in Section 4.3, “Using type, str, dir, and Other Built-In Functions” with regards to the name of the <code>urllib</code> module. Thanks to Rick for the correction.</p> <p>Added additional example in Section 16.2, “Finding the path” to show how to run unit tests in the current working directory, instead of the current directory.</p> <p>Modified explanation of how to derive a negative list index from a positive list index in Example 3.7, “Negative List Indices”. Thanks to Kevin for the suggestion.</p> <p>Updated links on home page for downloading latest version of Python.</p> <p>Added link on home page to Bruce Eckel’s preliminary draft of <i>Thinking in Python</i> (<a href="http://www.mindview.net/Books/TIPython">http://www.mindview.net/Books/TIPython</a>), a book about Python.</p>	
Revision 3.8	2001-11-18
<p>Added Section 16.2, “Finding the path”.</p> <p>Added Section 16.3, “Filtering lists revisited”.</p> <p>Added Section 16.4, “Mapping lists revisited”.</p> <p>Added Section 16.5, “Data-centric programming”.</p> <p>Expanded sample output in Section 16.1, “Diving in”.</p> <p>Finished Section 9.3, “Parsing XML”.</p>	
Revision 3.7	2001-09-30
<p>Added Section 9.2, “Packages”.</p>	

<p>Added Section 9.3, “Parsing XML”.</p> <p>Cleaned up introductory paragraph in Section 9.1, “Diving in”. Thanks to Matt for this suggestion.</p> <p>Added Java tip in Section 5.2, “Importing Modules Using from module import”. Thanks to Ori for this suggestion.</p> <p>Fixed mistake in Section 4.8, “Putting It All Together” where I implied that you could not use <code>None</code> to compare to a null value in Python.</p> <p>Clarified in Section 3.2, “Introducing Lists” where I said that <code>None</code> was equivalent to <code>False</code>. The result is the same, but <code>None</code> is faster because it doesn’t create a new object.</p> <p>Fixed mistake in Section 3.2, “Introducing Lists” where I said that <code>None</code> was equivalent to <code>False</code>. In fact, it’s equivalent to <code>False</code> since it doesn’t create a new object.</p> <p>Fixed typographical laziness in Chapter 2, <i>Your First Python Program</i>; when I was writing it, I had not yet standardized on putting <code>None</code> on a new line.</p> <p>Fixed mistake in Section 2.2, “Declaring Functions” where I said that statically typed languages always use explicit variable + data type.</p> <p>Added link to Spanish translation (<a href="http://es.diveintopython.org/">http://es.diveintopython.org/</a>).</p>	
Revision 3.6.4	2001-09-06
<p>Added code in <code>None</code> to handle non-HTML entity references, and added a note about it in Section 8.4, “Introducing BaseHTMLProcessor”.</p> <p>Modified Example 8.11, “Introducing globals” to include <code>None</code> in the output.</p>	
Revision 3.6.3	2001-09-04
<p>Fixed typo in Section 9.1, “Diving in”.</p> <p>Added link to Korean translation (<a href="http://kr.diveintopython.org/html/index.htm">http://kr.diveintopython.org/html/index.htm</a>).</p>	
Revision 3.6.2	2001-08-31
Fixed typo in Section 13.6, “Testing for sanity” (the last requirement was listed twice).	
Revision 3.6	2001-08-31
<p>Finished Chapter 8, <i>HTML Processing</i> with Section 8.9, “Putting it all together” and Section 8.10, “Summary”.</p> <p>Added Section 15.4, “Postscript”.</p> <p>Started Chapter 9, <i>XML Processing</i> with Section 9.1, “Diving in”.</p> <p>Started Chapter 16, <i>Functional Programming</i> with Section 16.1, “Diving in”.</p> <p>Fixed long-standing bug in colorizing script that improperly colorized the examples in Chapter 8, <i>HTML Processing</i>.</p> <p>Added link to French translation (<a href="http://fr.diveintopython.org/toc.html">http://fr.diveintopython.org/toc.html</a>). They did the right thing and translated the source XML, so</p> <p>Upgraded to version 1.43 of the DocBook XSL stylesheets.</p> <p>Upgraded to version 6.43 of the SAXON XSLT processor from Michael Kay.</p> <p>Massive stylesheet changes, moving away from a table-based layout and towards more appropriate use of cascading style sheets. Un</p> <p>Moved to Ant (<a href="http://jakarta.apache.org/ant/">http://jakarta.apache.org/ant/</a>) to have better control over the build process, which is especially important now that I’</p> <p>Consolidated the available downloadable archives; previously, I had different files for each platform, because the .zip files that Pyth</p> <p>Now hosting the complete XML source, XSL stylesheets, and associated scripts and libraries on SourceForge. There’s also CVS acc</p> <p>Re-licensed the example code under the new-and-improved GPL-compatible Python 2.1.1 license (<a href="http://www.python.org/2.1.1/lice">http://www.python.org/2.1.1/lice</a></p>	
Revision 3.5	2001-06-26
<p>Added explanation of strong/weak/static/dynamic datatypes in Section 2.2, “Declaring Functions”.</p> <p>Added case-sensitivity example in Section 3.1, “Introducing Dictionaries”.</p> <p>Use <code>None</code> in Chapter 5, <i>Objects and Object-Oriented</i> to compensate for inferior operating systems whose files aren’t case-sensitive.</p> <p>Fixed indentation problems in code samples in PDF version.</p>	
Revision 3.4	2001-05-31
<p>Added Section 14.5, “roman.py, stage 5”.</p> <p>Added Section 15.1, “Handling bugs”.</p> <p>Added Section 15.2, “Handling changing requirements”.</p> <p>Added Section 15.3, “Refactoring”.</p> <p>Added Section 15.5, “Summary”.</p> <p>Fixed yet another stylesheet bug that was dropping nested <code>&lt;div&gt;</code> tags.</p>	
Revision 3.3	2001-05-24
<p>Added Section 13.2, “Diving in”.</p> <p>Added Section 13.3, “Introducing romantest.py”.</p> <p>Added Section 13.4, “Testing for success”.</p> <p>Added Section 13.5, “Testing for failure”.</p> <p>Added Section 13.6, “Testing for sanity”.</p>	

<p>Added Section 14.1, “roman.py, stage 1”.</p> <p>Added Section 14.2, “roman.py, stage 2”.</p> <p>Added Section 14.3, “roman.py, stage 3”.</p> <p>Added Section 14.4, “roman.py, stage 4”.</p> <p>Tweaked stylesheets in an endless quest for complete Netscape/Mozilla compatibility.</p>	
Revision 3.2	2001-05-03
<p>Added Section 8.8, “Introducing dialect.py”.</p> <p>Added Section 7.2, “Case Study: Street Addresses”.</p> <p>Fixed bug in <code>method</code> that would produce incorrect declarations (adding a space where it couldn’t be).</p> <p>Fixed bug in CSS (introduced in 2.9) where body background color was missing.</p>	
Revision 3.1	2001-04-18
<p>Added code in <code>no</code> to handle declarations, now that Python 2.1 supports them.</p> <p>Added note about nested scopes in Section 8.5, “locals and globals”.</p> <p>Fixed obscure bug in Example 8.1, “BaseHTMLProcessor.py” where attribute values with character entities would not be properly c</p> <p>Now recommending (but not requiring) Python 2.1, due to its support of declarations in <code>h</code></p> <p>Updated download links on the home page (<a href="http://diveintopython.org/">http://diveintopython.org/</a>) to point to Python 2.1, where available.</p> <p>Moved to versioned filenames, to help people who redistribute the book.</p>	
Revision 3.0	2001-04-16
<p>Fixed minor bug in code listing in Chapter 8, <i>HTML Processing</i>.</p> <p>Added link to Chinese translation on home page (<a href="http://diveintopython.org/">http://diveintopython.org/</a>).</p>	
Revision 2.9	2001-04-13
<p>Added Section 8.5, “locals and globals”.</p> <p>Added Section 8.6, “Dictionary-based string formatting”.</p> <p>Tightened code in Chapter 8, <i>HTML Processing</i>, specifically <code>to</code> to use fewer and simpler regular expressions.</p> <p>Fixed a stylesheet bug that was inserting blank pages between chapters in the PDF version.</p> <p>Fixed a script bug that was stripping the <code>from</code> from the home page (<a href="http://diveintopython.org/">http://diveintopython.org/</a>).</p> <p>Added link to Python Cookbook (<a href="http://www.activestate.com/ASPN/Python/Cookbook/">http://www.activestate.com/ASPN/Python/Cookbook/</a>), and added a few links to individual recipe</p> <p>Switched to Google (<a href="http://www.google.com/services/free.html">http://www.google.com/services/free.html</a>) for searching on <code>h</code></p> <p>Upgraded to version 1.36 of the DocBook XSL stylesheets, which was much more difficult than it sounds. There may still be linger</p>	
Revision 2.8	2001-03-26
<p>Added Section 8.3, “Extracting data from HTML documents”.</p> <p>Added Section 8.4, “Introducing BaseHTMLProcessor.py”.</p> <p>Added Section 8.7, “Quoting attribute values”.</p> <p>Tightened up code in Chapter 4, <i>The Power Of Introspection</i>, using the built-in function <code>instead</code> of manually checking types.</p> <p>Moved Section 5.2, “Importing Modules Using from module import” from Chapter 4, <i>The Power Of Introspection</i> to Chapter 5, <i>Ob</i></p> <p>Fixed typo in code example in Section 5.1, “Diving In” (added colon).</p> <p>Added several additional downloadable example scripts.</p> <p>Added Windows Help output format.</p>	
Revision 2.7	2001-03-16
<p>Added Section 8.2, “Introducing sgmlib.py”.</p> <p>Tightened up code in Chapter 8, <i>HTML Processing</i>.</p> <p>Changed code in Chapter 2, <i>Your First Python Program</i> to use <code>method</code> instead of <code>h</code></p> <p>Moved Section 3.4.2, “Assigning Multiple Values at Once” section to Chapter 2, <i>Your First Python Program</i>.</p> <p>Edited note about <code>string</code> method, and provided a link to the new entry in <i>The Whole Python FAQ</i> (<a href="http://www.python.org/doc/FAQ">http://www.python.org/doc/FAQ</a>)</p> <p>Rewrote Section 4.6, “The Peculiar Nature of and and or” to emphasize the fundamental nature of <code>and</code> and <code>and</code> and de-emphasize the <code>and</code> trick</p> <p>Reorganized language comparisons into <code>h</code>.</p>	
Revision 2.6	2001-02-28
<p>The PDF and Word versions now have colorized examples, an improved table of contents, and properly indented <code>h</code> and <code>h</code>.</p> <p>The Word version is now in native Word format, compatible with Word 97.</p>	

<p>The PDF and text versions now have fewer problems with improperly converted special characters (like trademark symbols and currency symbols).</p> <p>Added link to download Word version for UNIX, in case some twisted soul wants to import it into StarOffice or something.</p> <p>Fixed several <b>h</b> which were missing titles.</p> <p>Fixed stylesheets to work around bug in Internet Explorer 5 for Mac OS which caused colorized words in the examples to be displayed in all caps.</p> <p>Fixed archive corruption in Mac OS downloads.</p> <p>In first section of each chapter, added link to download examples. (My access logs show that people skim or skip the two pages where the examples are.)</p> <p>Tightened the home page (<a href="http://diveintopython.org/">http://diveintopython.org/</a>) and preface even more, in the hopes that someday someone will read them.</p> <p>Soon I hope to get back to actually writing this book instead of debugging it.</p>	
Revision 2.5	2001-02-23
<p>Added Section 6.4, "Using sys.modules".</p> <p>Added Section 6.5, "Working with Directories".</p> <p>Moved Example 6.17, "Splitting Pathnames" from Section 3.4.2, "Assigning Multiple Values at Once" to Section 6.5, "Working with Directories".</p> <p>Added Section 6.6, "Putting It All Together".</p> <p>Added Section 5.10, "Summary".</p> <p>Added Section 8.1, "Diving in".</p> <p>Fixed program listing in Example 6.10, "Iterating Through a Dictionary" which was missing a colon.</p>	
Revision 2.4.1	2001-02-12
<p>Changed newsgroup links to use "news:" protocol, now that <b>g</b> is defunct.</p> <p>Added file sizes to download links.</p>	
Revision 2.4	2001-02-12
<p>Added "further reading" links in most sections, and collated them in Appendix A, <i>Further reading</i>.</p> <p>Added URLs in parentheses next to external links in text version.</p>	
Revision 2.3	2001-02-09
<p>Rewrote some of the code in Chapter 5, <i>Objects and Object-Oriented Programming</i> to use class attributes and a better example of multi-variables.</p> <p>Reorganized Chapter 5, <i>Objects and Object-Oriented Programming</i> to put the class sections first.</p> <p>Added Section 5.8, "Introducing Class Attributes".</p> <p>Added Section 6.1, "Handling Exceptions".</p> <p>Added Section 6.2, "Working with File Objects".</p> <p>Merged the "review" section in Chapter 5, <i>Objects and Object-Oriented Programming</i> into Section 5.1, "Diving In".</p> <p>Colorized all program listings and examples.</p> <p>Fixed important error in Section 2.2, "Declaring Functions": functions that do not explicitly return a value return <b>None</b>. So you can assign a variable to the result of a function call and it will be <b>None</b>.</p> <p>Added minor clarifications to Section 2.3, "Documenting Functions", Section 2.4, "Everything Is an Object", and Section 3.4, "Declaring Functions".</p>	
Revision 2.2	2001-02-02
<p>Edited Section 4.4, "Getting Object References With getattr".</p> <p>Added titles to <b>h</b> tags, so they can have their cute little tooltips too.</p> <p>Changed the look of the revision history page.</p> <p>Fixed problem I introduced yesterday in my HTML post-processing script that was causing invalid HTML character references and broken links.</p> <p>Upgraded to version 1.29 of the DocBook XSL stylesheets.</p>	
Revision 2.1	2001-02-01
<p>Rewrote the example code of Chapter 4, <i>The Power Of Introspection</i> to use <b>isinstance</b> instead of <b>isinstance</b> and <b>isinstance</b> and rewrote explanatory text to match.</p> <p>Added example of list operators in Section 3.2, "Introducing Lists".</p> <p>Added links to relevant sections in the summary lists at the end of each chapter (Section 3.8, "Summary" and Section 4.9, "Summary").</p>	
Revision 2.0	2001-01-31
<p>Split Section 5.6, "Special Class Methods" into three sections, Section 5.5, "Exploring UserDict: A Wrapper Class", Section 5.6, "Special Class Methods", and Section 5.7, "Special Class Methods".</p> <p>Changed notes on garbage collection to point out that Python 2.0 and later can handle circular references without additional coding.</p> <p>Fixed UNIX downloads to include all relevant files.</p>	
Revision 1.9	2001-01-15
<p>Removed introduction to Chapter 2, <i>Your First Python Program</i>.</p> <p>Removed introduction to Chapter 4, <i>The Power Of Introspection</i>.</p>	

Removed introduction to Chapter 5, <i>Objects and Object-Orientation</i> . Edited text ruthlessly. I tend to ramble.	
Revision 1.8	2001-01-12
Added more examples to Section 3.4.2, “Assigning Multiple Values at Once”. Added Section 5.3, “Defining Classes”. Added Section 5.4, “Instantiating Classes”. Added Section 5.6, “Special Class Methods”. More minor stylesheet tweaks, including adding titles to <code>h</code> ags, which, if your browser is cool enough, will display a description of the	
Revision 1.71	2001-01-03
Made several modifications to stylesheets to improve browser compatibility.	
Revision 1.7	2001-01-02
Added introduction to Chapter 2, <i>Your First Python Program</i> . Added introduction to Chapter 4, <i>The Power Of Introspection</i> . Added review section to Chapter 5, <i>Objects and Object-Orientation</i> [later removed] Added Section 5.9, “Private Functions”. Added Section 6.3, “Iterating with for Loops”. Added Section 3.4.2, “Assigning Multiple Values at Once”. Wrote scripts to convert book to new output formats: one single HTML file, PDF, Microsoft Word 97, and plain text. Registered the <code>h</code> omain and moved the book there, along with links to download the book in all available output formats for offline use. Modified the XSL stylesheets to change the header and footer navigation that displays on each page. The top of each page is branded	
Revision 1.6	2000-12-11
Added Section 4.8, “Putting It All Together”. Finished Chapter 4, <i>The Power Of Introspection</i> with Section 4.9, “Summary”. Started Chapter 5, <i>Objects and Object-Orientation</i> with Section 5.1, “Diving In”.	
Revision 1.5	2000-11-22
Added Section 4.6, “The Peculiar Nature of and and or”. Added Section 4.7, “Using lambda Functions”. Added appendix that lists section abstracts. Added appendix that lists tips. Added appendix that lists examples. Added appendix that lists revision history. Expanded example of mapping lists in Section 3.6, “Mapping Lists”. Encapsulated several more common phrases into entities. Upgraded to version 1.25 of the DocBook XSL stylesheets.	
Revision 1.4	2000-11-14
Added Section 4.5, “Filtering Lists”. Added <code>h</code> ocumentation to Section 4.3, “Using type, str, dir, and Other Built-In Functions”. Added <code>h</code> example in Section 3.3, “Introducing Tuples”. Added additional note about <code>h</code> rick under MacPython. Switched to the SAXON XSLT processor from Michael Kay. Upgraded to version 1.24 of the DocBook XSL stylesheets. Added db-html processing instructions with explicit filenames of each chapter and section, to allow deep links to content even if I a Made several common phrases into entities for easier reuse. Changed several <code>h</code> ags to <code>h</code>	
Revision 1.3	2000-11-09
Added section on dynamic code execution. Added links to relevant section/example wherever I refer to previously covered concepts. Expanded introduction of chapter 2 to explain what the function actually does. Explicitly placed example code under the GNU General Public License and added appendix to display license. [Note 8/16/2001: co	

Changed links to licenses to use flags, now that I know how to use them.	
Revision 1.2	2000-11-06
Added first four sections of chapter 2. Tightened up preface even more, and added link to Mac OS version of Python. Filled out examples in "Mapping lists" and "Joining strings" to show logical progression. Added output in chapter 1 summary.	
Revision 1.1	2000-10-31
Finished chapter 1 with sections on mapping and joining, and a chapter summary. Toned down the preface, added links to introductions for non-programmers. Fixed several typos.	
Revision 1.0	2000-10-30
Initial publication.	



## Appendix F. About the book

This book was written in DocBook XML (<http://www.oasis-open.org/docbook/>) using Emacs (<http://www.gnu.org/software/emacs/>), and con

If you're interested in learning more about DocBook for technical writing, you can download the XML source (<http://diveintopython.org/dow>

# Appendix G. GNU Free Documentation License

Version 1.1, March 2000

Copyright (C) 2000 Free Software Foundation, Inc. 59 Temple Place, Suite 330, Boston, MA 02111-1307 USA Everyone is permitted to copy and distribute verbatim copies of this license document, provided the copyright notice and this notice are preserved on all copies.

## G.0. Preamble

The purpose of this License is to make a manual, textbook, or other written document "free" in the sense of freedom: to assure everyone the e

This License is a kind of "copyleft", which means that derivative works of the document must themselves be free in the same sense. It compl

We have designed this License in order to use it for manuals for free software, because free software needs free documentation: a free progra

## G.1. Applicability and definitions

This License applies to any manual or other work that contains a notice placed by the copyright holder saying it can be distributed under the t

A "Modified Version" of the Document means any work containing the Document or a portion of it, either copied verbatim, or with modifica

A "Secondary Section" is a named appendix or a front-matter section of the Document that deals exclusively with the relationship of the publi

The "Invariant Sections" are certain Secondary Sections whose titles are designated, as being those of Invariant Sections, in the notice that sa

The "Cover Texts" are certain short passages of text that are listed, as Front-Cover Texts or Back-Cover Texts, in the notice that says that the

A "Transparent" copy of the Document means a machine-readable copy, represented in a format whose specification is available to the gener

Examples of suitable formats for Transparent copies include plain ASCII without markup, Texinfo input format, LaTeX input format, SGML

The "Title Page" means, for a printed book, the title page itself, plus such following pages as are needed to hold, legibly, the material this Lic

## G.2. Verbatim copying

You may copy and distribute the Document in any medium, either commercially or noncommercially, provided that this License, the copyrig

You may also lend copies, under the same conditions stated above, and you may publicly display copies.

## G.3. Copying in quantity

If you publish printed copies of the Document numbering more than 100, and the Document's license notice requires Cover Texts, you must

If the required texts for either cover are too voluminous to fit legibly, you should put the first ones listed (as many as fit reasonably) on the ac

If you publish or distribute Opaque copies of the Document numbering more than 100, you must either include a machine-readable Transpare

It is requested, but not required, that you contact the authors of the Document well before redistributing any large number of copies, to give th

## G.4. Modifications

You may copy and distribute a Modified Version of the Document under the conditions of sections 2 and 3 above, provided that you release t

Use in the Title Page (and on the covers, if any) a title distinct from that of the Document, and from those of previous versions (whic

List on the Title Page, as authors, one or more persons or entities responsible for authorship of the modifications in the Modified Ve

State on the Title page the name of the publisher of the Modified Version, as the publisher.3.

Preserve all the copyright notices of the Document.4.

Add an appropriate copyright notice for your modifications adjacent to the other copyright notices.5.

Include, immediately after the copyright notices, a license notice giving the public permission to use the Modified Version under the

Preserve in that license notice the full lists of Invariant Sections and required Cover Texts given in the Document's license notice.7.

Include an unaltered copy of this License.8.

Preserve the section entitled "History", and its title, and add to it an item stating at least the title, year, new authors, and publisher of

Preserve the network location, if any, given in the Document for public access to a Transparent copy of the Document, and likewise

In any section entitled "Acknowledgements" or "Dedications", preserve the section's title, and preserve in the section all the substan

Preserve all the Invariant Sections of the Document, unaltered in their text and in their titles. Section numbers or the equivalent are n

Delete any section entitled "Endorsements". Such a section may not be included in the Modified Version.13.

Do not retitle any existing section as "Endorsements" or to conflict in title with any Invariant Section.14.

If the Modified Version includes new front-matter sections or appendices that qualify as Secondary Sections and contain no material copied f

You may add a section entitled "Endorsements", provided it contains nothing but endorsements of your Modified Version by various parties-

You may add a passage of up to five words as a Front-Cover Text, and a passage of up to 25 words as a Back-Cover Text, to the end of the li

The author(s) and publisher(s) of the Document do not by this License give permission to use their names for publicity for or to assert or imp

## G.5. Combining documents

You may combine the Document with other documents released under this License, under the terms defined in section 4 above for modified v

The combined work need only contain one copy of this License, and multiple identical Invariant Sections may be replaced with a single copy

In the combination, you must combine any sections entitled "History" in the various original documents, forming one section entitled "Histor

## G.6. Collections of documents

You may make a collection consisting of the Document and other documents released under this License, and replace the individual copies o

You may extract a single document from such a collection, and distribute it individually under this License, provided you insert a copy of thi

## G.7. Aggregation with independent works

A compilation of the Document or its derivatives with other separate and independent documents or works, in or on a volume of a storage or

If the Cover Text requirement of section 3 is applicable to these copies of the Document, then if the Document is less than one quarter of the

## G.8. Translation

Translation is considered a kind of modification, so you may distribute translations of the Document under the terms of section 4. Replacing

## G.9. Termination

You may not copy, modify, sublicense, or distribute the Document except as expressly provided for under this License. Any other attempt to

## G.10. Future revisions of this license

The Free Software Foundation may publish new, revised versions of the GNU Free Documentation License from time to time. Such new ver

Each version of the License is given a distinguishing version number. If the Document specifies that a particular numbered version of this Lic

## G.11. How to use this License for your documents

To use this License in a document you have written, include a copy of the License in the document and put the following copyright and license notice in the document:

Copyright (c) YEAR YOUR NAME. Permission is granted to copy, distribute and/or modify this document under the terms of the C

If you have no Invariant Sections, write "with no Invariant Sections" instead of saying which ones are invariant. If you have no Front-Cover Texts, write "with no Front-Cover Texts" instead of saying which ones are invariant. If you have no Back-Cover Texts, write "with no Back-Cover Texts" instead of saying which ones are invariant. If you have no Invariant License, write "with no Invariant License" instead of saying which ones are invariant. If you have no Invariant Sections, write "with no Invariant Sections" instead of saying which ones are invariant. If you have no Front-Cover Texts, write "with no Front-Cover Texts" instead of saying which ones are invariant. If you have no Back-Cover Texts, write "with no Back-Cover Texts" instead of saying which ones are invariant. If you have no Invariant License, write "with no Invariant License" instead of saying which ones are invariant.

If your document contains nontrivial examples of program code, we recommend releasing these examples in parallel under your choice of free

# Appendix H. Python license

## H.A. History of the software

Python was created in the early 1990s by Guido van Rossum at Stichting Mathematisch Centrum (CWI) in the Netherlands as a successor of

Following the release of Python 1.6, and after Guido van Rossum left CNRI to work with commercial software developers, it became clear th

After Python 2.0 was released by BeOpen.com, Guido van Rossum and the other PythonLabs developers joined Digital Creations. All intelle

Thanks to the many outside volunteers who have worked under Guido's direction to make these releases possible.

## H.B. Terms and conditions for accessing or otherwise using Python

### H.B.1. PSF license agreement

This LICENSE AGREEMENT is between the Python Software Foundation ("PSF"), and the Individual or Organization ("Licensee")  
Subject to the terms and conditions of this License Agreement, PSF hereby grants Licensee a nonexclusive, royalty-free, world-wide  
In the event Licensee prepares a derivative work that is based on or incorporates Python 2.1.1 or any part thereof, and wants to make  
PSF is making Python 2.1.1 available to Licensee on an "AS IS" basis. PSF MAKES NO REPRESENTATIONS OR WARRANTIES  
PSF SHALL NOT BE LIABLE TO LICENSEE OR ANY OTHER USERS OF PYTHON 2.1.1 FOR ANY INCIDENTAL, SPECIAL, OR CONSEQUENTIAL DAMAGES OR LOSS.  
This License Agreement will automatically terminate upon a material breach of its terms and conditions.  
Nothing in this License Agreement shall be deemed to create any relationship of agency, partnership, or joint venture between PSF and  
By copying, installing or otherwise using Python 2.1.1, Licensee agrees to be bound by the terms and conditions of this License Agreement.

### H.B.2. BeOpen Python open source license agreement version 1

This LICENSE AGREEMENT is between BeOpen.com ("BeOpen"), having an office at 160 Saratoga Avenue, Santa Clara, CA 95051  
Subject to the terms and conditions of this BeOpen Python License Agreement, BeOpen hereby grants Licensee a non-exclusive, royalty-free, world-wide  
BeOpen is making the Software available to Licensee on an "AS IS" basis. BEOPEN MAKES NO REPRESENTATIONS OR WARRANTIES  
BEOPEN SHALL NOT BE LIABLE TO LICENSEE OR ANY OTHER USERS OF THE SOFTWARE FOR ANY INCIDENTAL, SPECIAL, OR CONSEQUENTIAL DAMAGES OR LOSS.  
This License Agreement will automatically terminate upon a material breach of its terms and conditions.  
This License Agreement shall be governed by and interpreted in all respects by the law of the State of California, excluding conflict of laws provisions.  
By copying, installing or otherwise using the software, Licensee agrees to be bound by the terms and conditions of this License Agreement.

### H.B.3. CNRI open source GPL-compatible license agreement

This LICENSE AGREEMENT is between the Corporation for National Research Initiatives, having an office at 1895 Preston White Drive  
Subject to the terms and conditions of this License Agreement, CNRI hereby grants Licensee a nonexclusive, royalty-free, world-wide  
In the event Licensee prepares a derivative work that is based on or incorporates Python 1.6.1 or any part thereof, and wants to make  
CNRI is making Python 1.6.1 available to Licensee on an "AS IS" basis. CNRI MAKES NO REPRESENTATIONS OR WARRANTIES  
CNRI SHALL NOT BE LIABLE TO LICENSEE OR ANY OTHER USERS OF PYTHON 1.6.1 FOR ANY INCIDENTAL, SPECIAL, OR CONSEQUENTIAL DAMAGES OR LOSS.  
This License Agreement will automatically terminate upon a material breach of its terms and conditions.  
This License Agreement shall be governed by the federal intellectual property law of the United States, including without limitation the patent laws.  
By clicking on the "ACCEPT" button where indicated, or by copying, installing or otherwise using Python 1.6.1, Licensee agrees to be bound by the terms and conditions of this License Agreement.

### H.B.4. CWI permissions statement and disclaimer

Copyright (c) 1991 - 1995, Stichting Mathematisch Centrum Amsterdam, The Netherlands. All rights reserved.

Permission to use, copy, modify, and distribute this software and its documentation for any purpose and without fee is hereby granted, provided that

STICHTING MATHEMATISCH CENTRUM DISCLAIMS ALL WARRANTIES WITH REGARD TO THIS SOFTWARE, INCLUDING ANY IMPLIED WARRANTY OF MERCHANTABILITY.