# Table of Contents

# Table of Contents

# Table of Contents

# Table of Contents

# Table of Contents

# 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 or any later version published by the Free Software Foundation; with no Invariant Sections, no Front–Cover Texts, and no Back–Cover Texts. A copy of the license is included in Appendix G, *GNU Free Documentation License*.

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 Python Software Foundation. A copy of the license is included in Appendix H, *Python license*.

# 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 in the default installation. Mac OS X 10.2 and later includes a command–line version of Python, although you'll probably want to install a version that includes a more Mac–like graphical interface.

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 free UNIX–compatible systems like Linux. There are also versions that run on Sun Solaris, AS/400, Amiga, OS/2, BeOS, and a plethora of other platforms you've probably never even heard of.

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 Windows and later deploy them 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–aware code editor, plus some Windows extensions for Python that allow complete access to Windows–specific services, APIs, and the Windows Registry.

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 specific reason not to. One such reason might be that ActiveState is generally several months behind in updating their ActivePython installer when new version of Python are released. If you absolutely need the latest version of Python and ActivePython is still a version behind as you read this, you'll want to use the second option for installing Python on Windows.

The second option is the "official" Python installer, distributed by the people who develop Python itself. It is freely downloadable and open source, and it is always current with the latest version of Python.

**Procedure 1.1. Option 1Lent k1 Tf 0 eo open sourcelatOn, which incoclul itself.  a9n procedurstill a versionel itcMpen**

3. Double−click the installer, `ActivePython-2.2.2-224-win32-ix86.msi`.
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 can't spare the 14MB.
6. After the installation is complete, close the installer and choose Start−>Programs−>ActiveState ActivePython 2.2−>PythonWin IDE. You'll see something like the following:

```
PythonWin 2.2.2 (#37, Nov 26 2002, 10:24:37) [MSC 32 bit (Intel)] on win32.
Portions Copyright 1994-2001 Mark Hammond (mhammond@skippinet.com.au) −
see 'Help/About PythonWin' for further copyright information.
>>>
```

**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 listed, then downloading the `.exe` installer.
2. Double−click the installer, `Python-2.xxx.yyy.exe`. 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 (`Tools/`), and/or the test suite (`Lib/test/`).
5. If you do not have administrative rights on your machine, you can select Advanced Options, then choose Non−Admin Install. This just affects where Registry entries and Start menu shortcuts are created.
6. After the installation is complete, close the installer and select Start−>Programs−>Python 2.3−>IDLE (Python GUI). You'll see something like the following:

```
Python 2.3.2 (#49, Oct  2 2003, 20:02:00) [MSC v.1200 32 bit (Intel)] on win32
Type "copyright", "credits" or "license()" for more information.

    ****************************************************************
    Personal firewall software may warn about the connection IDLE
    makes to its subprocess using this computer's internal loopback
    interface.  This connection is not visible on any external
    interface and no data is sent to or received from the Internet.
    ****************************************************************

IDLE 1.0
>>>
```

# 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 use this version for the first third of the book. However, the preinstalled version does not come with an XML parser, so when you get to the XML chapter, you'll need to install the full version.

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 `/Applications` folder.

2. Open the `Utilities` folder.

3. Double–click `Terminal` to open a terminal window and get to a command line.

4. Type **python** at the command prompt.

Try it out:

```
Welcome to Darwin!
[localhost:~] you% python
Python 2.2 (#1, 07/14/02, 23:25:09)
[GCC Apple cpp-precomp 6.14] on darwin
Type "help", "copyright", "credits", or "license" for more information.
>>> [press Ctrl+D to get back to the command prompt]
[localhost:~] you%
```

**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 `MacPython-OSX` disk image from http://homepages.cwi.nl/~jack/macpython/download.html. If your browser has not already done so, double–click `MacPython-OSX-2.3-1.dmg` to mount the disk

## 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 `MacPython23full.bin` file from http://homepages.cwi.nl/~jack/macpython/download.html.
2. If your browser does not decompress the file automatically, double–click `MacPython23full.bin` to decompress the file with Stuffit Expander.
3. Double–click the installer, `MacPython23full`.
4. Step through the installer program.
5. AFter installation is complete, close the installer and open the `/Applications` folder.
6. Open the `MacPython-OS9 2.3` folder.
7. Double–click `Python IDE` 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 Window–>Python Interactive (**Cmd–0**). You'll see a screen like this:

```
Python 2.3 (#2, Jul 30 2003, 11:45:28)
[GCC 3.1 20020420 (prerelease)]
Type "copyright", "credits" or "license" for more information.
MacPython IDE 1.0.1
>>>
```

## 1.5. Python on RedHat Linux

Installing under UNIX–compatible operating systems such as Linux is easy if you're willing to install a binary package. Pre–built binary packages are available for most popular Linux distributions. Or you can always compile from source.

Download the latest Python RPM by going to http://www.python.org/ftp/python/ and selecting the highest version number listed, then selecting the `rpms/` dir9.2lecta .R:

❶

❷

```
Python 2.3 (#1, Sep 12 2003, 10:53:56)
[GCC 3.2.2 20030222 (Red Hat Linux 3.2.2-5)] on linux2
Type "help", "copyright", "credits", or "license" for more information.
>>> [press Ctrl+D to exit]
[root@localhost root]# which python2.3 ❸
/usr/bin/python2.3
```

❶    Whoops! Just typing **python** gives you the older version of Python –– the one that was installed by default. That's not the one you want.

❷    At the time of this writing, the newest version is called **python2.3**. You'll probably want to change the path on the first line of the sample scripts to point to the newer version.

❸    This is the complete path of the newer version of Python that you just installed. Use this on the #! line (the first line of each script) to ensure that scripts are running under the latest version of Python, and be sure to type **python2.3** to get into the interactive shell.

## 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.

**Example 1.3. Installing on Debian GNU/Linux**

```
localhost:~$ su -
Password: [enter your root password]
localhost:~# apt-get install python
Reading Package Lists... Done
Building Dependency Tree... Done
The following extra packages will be installed:
  python2.3
Suggested packages:
  python-tk python2.3-doc
The following NEW packages will be installed:
  python python2.3
0 upgraded, 2 newly installed, 0 to remove and 3 not upgraded.
Need to get 0B/2880kB of archives.
After unpacking 9351kB of additional disk space will be used.
Do you want to continue? [Y/n] Y
Selecting previously deselected package python2.3.
(Reading database ... 22848 files and directories currently installed.)
Unpacking python2.3 (from .../python2.3_2.3.1-1_i386.deb) ...
Selecting previously deselected package python.
Unpacking python (from .../python_2.3.1-1_all.deb) ...
Setting up python (2.3.1-1) ...
Setting up python2.3 (2.3.1-1) ...
Compiling python modules in /usr/lib/python2.3 ...
Compiling optimized python modules in /usr/lib/python2.3 ...
localhost:~# exit
logout
localhost:~$ python
Python 2.3.1 (#2, Sep 24 2003, 11:39:14)
[GCC 3.3.2 20030908 (Debian prerelease)] on linux2
Type "help", "copyright", "credits" or "license" for more information.
>>> [press Ctrl+D to exit]
```

## 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 number listed, download the `.tgz` file), and then do the usual **configure**, **make**, **make**

**install** dance.


## Example 1.4. Installing from source

```
localhost:~$ su -
Password: [enter your root password]
localhost:~# wget http://www.python.org/ftp/python/2.3/Python-2.3.tgz
Resolving www.python.org... done.
Connecting to www.python.org[194.109.137.226]:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: 8,436,880 [application/x-tar]
...
localhost:~# tar xfz Python-2.3.tgz
localhost:~# cd Python-2.3
localhost:~/Python-2.3# ./configure
checking MACHDEP... linux2
checking EXTRAPLATDIR...
checking for --without-gcc... no
...
localhost:~/Python-2.3# make
gcc -pthread -c -fno-strict-aliasing -DNDEBUG -g -O3 -Wall -Wstrict-prototypes
-I. -I./Include  -DPy_BUILD_CORE -o Modules/python.o Modules/python.c
gcc -pthread -c -fno-strict-aliasing -DNDEBUG -g -O3 -Wall -Wstrict-prototypes
-I. -I./Include  -DPy_BUILD_CORE -o Parser/acceler.o Parser/acceler.c
gcc -pthread -c -fno-strict-aliasing -DNDEBUG -g -O3 -Wall -Wstrict-prototypes
-I. -I./Include  -DPy_BUILD_CORE -o Parser/grammar1.o Parser/grammar1.c
...
localhost:~/Python-2.3# make install
/usr/bin/install -c python /usr/local/bin/python2.3
...
localhost:~/Python-2.3# exit
logout
localhost:~$ which python
/usr/local/bin/python
localhost:~$ python
Python 2.3.1 (#2, Sep 24 2003, 11:39:14)
[GCC 3.3.2 20030908 (Debian prerelease)] on linux2
```

❶

```
    Returns string."""
```

Triple quotes signify a multi–line string. Everything between the start and end quotes is part of a single string, including carriage returns and other quote characters. You can use them anywhere, but you'll see them most often used when defining a `doc string`.

Triple quotes are also an easy way to define a string with both single and double quotes, like `qq/.../` in Perl. Everything between the triple quotes is the function's `doc string`, which documents what the function does. A `doc string`, if it exists, must be the first thing defined in a function (that is, the first thing after the colon). You don't technically need to give your function a `doc string`, but you always should. I know you've heard this in every programming class you've ever taken, but Python gives you an added incentive: the `doc string` is available at runtime as an attribute of the function.

Many Python IDEs use the `doc string` to provide context–sensitive documentation, so that when you type a function name, its `doc string` appears as a tooltip. This can be incredibly helpful, but it's only as good as the `doc strings` you write.

**Further Reading on Documenting Functions**

- PEP 257 (http://www.python.org/peps/pep–0257.html) defines `doc string` conventions.
- *Python Style Guide* (http://www.python.org/doc/essays/styleguide.html) discusses how to write a good `doc string`.
- *Python Tutorial* (http://www.python.org/doc/current/tut/tut.html) discusses conventions for spacing in `doc strings` (http://www.python.org/doc/current/tut/node6.html#SECTION006750000000000000000).

## 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 `buildConnectionString` Function's `doc string`**

```
>>> import odbchelper                                      ❶
>>> params = {"server":"mpilgrim", "database":"master", "uid":"sa", "pwd":"secret"}
>>> print odbchelper.buildConnectionString(params) ❷
server=mpilgrim;uid=sa;database=master;pwd=secret
>>> print odbchelper.buildConnectionString.__doc__ ❸
```

❶

❷     When you want to use functions defined in imported modules, you need to include the module name. So you can't just say `buildConnectionString`; it must be `odbchelper.buildConnectionString`. If you've used classes in Java, this should feel vaguely familiar.

❸     Instead of calling the function as you would expect to, you asked for one of the function's attributes, `__doc__`.

`import` in Python is like `require` 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. Specifically, it looks in all the directories defined in `sys.path`. This is just a list, and you can easily view it or modify it with standard list methods. (You'll learn more about lists later in this chapter.)

**Example 2.4. Import Search Path**

```
>>> import sys                    ❶
>>> sys.path                      ❷
['', '/usr/local/lib/python2.2', '/usr/local/lib/python2.2/plat-linux2',
'/usr/local/lib/python2.2/lib-dynload', '/usr/local/lib/python2.2/site-packages',
'/usr/local/lib/python2.2/site-packages/PIL', '/usr/local/lib/python2.2/site-packages/piddle']
>>> sys                           ❸
<module 'sys' (built-in)>
>>> sys.path.append('/my/new/path')  ❹
```

❶     Importing the `sys` module makes all of its functions and attributes available.

❷     `sys.path` is a list of directory names that constitute the current search path. (Yours will look different, depending on your operating system, what version of Python you're running, and where it was originally installed.) Python will look through these directories (in this order) for a `.py` file matching the module name you're trying to import.

❸     Actually, I lied; the truth is more complicated than that, because not all modules are stored as `.py` files. Some, like the `sys` module, are "built–in modules"; they are actually baked right into Python itself. Built–in modules behave just like regular modules, but their Python source code is not available, because they are not written in Python! (The `sys` module is written in C.)

❹     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 as well, whenever you try to import a module. The effect lasts as long as Python is running. (You'll talk more about `append` and other list methods in Chapter 3.)

## 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 `doc string` defined in the function's source code. The `sys` module is an object which has (among other things) an attribute called `path`. And so forth.

Still, this begs the question. What is an object? Different programming languages define "object" in different ways. In some, it means that *all* objects *must* have attributes and methods; in others, it means that all objects are subclassable. In Python, the definition is looser; some objects have neither attributes nor methods (more on this in Chapter 3), and not all objects are subclassable (more on this in Chapter 5). But everything is an object in the sense that it can be assigned to a variable or passed as an argument to a function (more in this in Chapter 4).

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. Lists are objects. Functions are objects. Even modules are objects.

**Further Reading on Objects**

*Python Reference Manual*

❶
❷
❸

❹

❶
❷

❸

❹

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 all Python programs look similar, since indentation is a language requirement and not a matter of style. This makes it easier to read and understand other people's Python 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 and curly braces to separate code blocks.

**Further Reading on Code Indentation**

- *Python Reference Manual* (http://www.python.org/doc/current/ref/) discusses cross−platform indentation issues and shows various indentation errors (http://www.python.org/doc/current/ref/indentation.html).
- *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 that uses the if __name__ trick.

```
if __name__ == "__main__":
```

Some quick observations before you get to the good stuff. First, parentheses are not required around the if expression. Second, the if statement ends with a colon, and is followed by indented code.

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 assigning the value you thought you were comparing.

So why is this particular if statement a trick? Modules are objects, and all modules have a built−in attribute __name__. A module's __name__ depends on how you're using the module. If you import the module, then __name__ is the module's filename, without a directory path or file extension. But you can also run the module directly as a standalone program, in which case __name__ will be a special default value, __main__.

```
>>> import odbchelper
>>> odbchelper.__name__
'odbchelper'
```

Knowing this, you can design a test suite for your module within the module itself by putting it in this if statement. When you run the module directly, __name__ is __main__, so the test suite executes. When you import the module, __name__ is something else, so the test suite is ignored. This makes it easier to develop and debug new modules before integrating them into a larger program.

On MacPython, there is an additional step to make the if __name__ trick work. Pop up the module's options menu by clicking the black triangle in the upper−right corner of the window, and make sure Run as __main__ is checked.

**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/import.html).

# 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, tuples, and lists (oh my!). If you're a Perl hacker, you can probably skim the bits about dictionaries and lists, but you should still pay attention to tuples.

## 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, and Python keeps track of the datatype internally.

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

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

### 3.1.1. Defining Dictionaries

**Example 3.1. Defining a Dictionary**

```
>>> d = {"server":"mpilgrim", "database":"master"}    ❶
>>> d
{'server': 'mpilgrim', 'database': 'master'}
>>> d["server"]                                        ❷
'mpilgrim'
>>> d["database"]                                      ❸
'master'
>>> d["mpilgrim"]                                      ❹
Traceback (innermost last):
  File "<interactive input>", line 1, in ?
KeyError: mpilgrim
```

❶ 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 of elements is enclosed in curly braces.

❷ `'server'` is a key, and its associated value, referenced by d["server"], is 'mpilgrim'.

❸ `'database'` is a key, and its associated value, referenced by d["database"], is 'master'.

❹ You can get values by key, but you can't get keys by value. So d["server"] is 'mpilgrim', but d["mpilgrim"] raises an exception, because 'mpilgrim' is not a key.

### 3.1.2. Modifying Dictionaries

**Example 3.2. Modifying a Dictionary**

```
>>> d
{'server': 'mpilgrim', 'database': 'master'}
>>> d["database"] = "pubs"    ❶
>>> d
{'server': 'mpilgrim', 'database': 'pubs'}
>>> d["uid"] = "sa"           ❷
>>> d
{'server': 'mpilgrim', 'uid': 'sa', 'database': 'pubs'}
```

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 when you think you are adding new values but are actually just modifying the same value over and over because your key isn't changing the way you think it is.)

Note that the new element (key

❶

❷

❶
❷

❶

❷

❶

❷

### 3.1.3. Deleting Items From Dictionaries

**Example 3.5. Deleting Items from a Dictionary**

```
>>> d
{'server': 'mpilgrim', 'uid': 'sa', 'database': 'master',
42: 'douglas', 'retrycount': 3}
>>> del d[42]  ❶
>>> d
{'server': 'mpilgrim', 'uid': 'sa', 'database': 'master', 'retrycount': 3}
>>> d.clear()  ❷
>>> d
{}
```

❶  `del` lets you delete individual items from a dictionary by key.

❷  `clear` 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 dictionaries to model sparse matrices (http://www.ibiblio.org/obp/thinkCSpy/chap10.htm).
- Python Knowledge Base (http://www.faqts.com/knowledge−base/index.phtml/fid/199/) has a lot of example code using dictionaries (http://www.faqts.com/knowledge−base/index.phtml/fid/541).
- Python Cookbook (http://www.activestate.com/ASPN/Python/Cookbook/) discusses how to sort the values of a dictionary by key (http://www.activestate.com/ASPN/Python/Cookbook/Recipe/52306).
- *Python Library Reference* (http://www.python.org/doc/current/lib/) summarizes all the dictionary methods (http://www.python.org/doc/current/lib/typesmapping.html).

## 3.2. Introducing Lists

Lists are Python's workhorse datatype. If your only experience with lists is arrays in Visual Basic or (God forbid) thbrary"icual

❶

❷
❸

❶
❷

❶

❷

❶
❷
❸

❶

❸ Note the symmetry here. In this five–element list, `li[:3]` returns the first 3 elements, and `li[3:]` returns the last two elements. In fact, `li[:n]` will always return the first `n` elements, and `li[n:]` will return the rest, regardless of the length of the list.

❹ If both slice indices are left out, all elements of the list are included. But this is not the same as the original `li` list; it is a new list that happens to have all the same elements. `li[:]` is shorthand for making a complete copy of a list.

## 3.2.2. Adding Elements to Lists

**Example 3.10. Adding Elements to a List**

```
>>> li
['a', 'b', 'mpilgrim', 'z', 'example']
>>> li.append("new")                      ❶
>>> li
['a', 'b', 'mpilgrim', 'z', 'example', 'new']
>>> li.insert(2, "new")                    ❷
>>> li
['a', 'b', 'new', 'mpilgrim', 'z', 'example', 'new']
>>> li.extend(["two", "elements"])   ❸
>>> li
['a', 'b', 'new', 'mpilgrim', 'z', 'example', 'new', 'two', 'elements']
```

❶ `append` adds a single element to the end of the list.

❷ `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 list elements do not need to be unique; there are now two separate elements with the value `'new'`, `li[2]` and `li[6]`.

❸ `extend` concatenates lists. Note that you do not call `extend` with multiple arguments; you call it with one argument, a list. In this case, that list has two elements.

**Example 3.11. The Difference between `extend` and `append`**

```
>>> li = ['a', 'b', 'c']
>>> li.extend(['d', 'e', 'f'])  ❶
>>> li
['a', 'b', 'c', 'd', 'e', 'f']
>>> len(li)                        ❷
6
>>> li[-1]
'f'
>>> li = ['a', 'b', 'c']
>>> li.append(['d', 'e', 'f'])  ❸
>>> li
['a', 'b', 'c', ['d', 'e', 'f']]
>>> len(li)                        ❹
4
>>> li[-1]
['d', 'e', 'f']
```

❶ Lists have two methods, `extend` and `append`, that look like they do the same thing, but are in fact completely different. `extend` takes a single argument, which is always a list, and adds each of the elements of that list to the original list.

❷ 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.

❸

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 with a single argument, which is a list of three elements.

❹ Now the original list, which started as a list of three elements, contains four elements. Why four? Because the last element that you just appended *is itself a list*. Lists can contain any type of data, including other lists. That may be what you want, or maybe not. Don't use `append` if you mean `extend`.

## 3.2.3. Searching Lists

**Example 3.12. Searching a List**

```
>>> li
['a', 'b', 'new', 'mpilgrim', 'z', 'example', 'new', 'two', 'elements']
>>> li.index("example")  ❶
5
>>> li.index("new")      ❷
2
>>> li.index("c")        ❸
Traceback (innermost last):
  File "<interactive input>", line 1, in ?
ValueError: list.index(x): x not in list
>>> "c" in li            ❹
False
```

❶ `index` finds the first occurrence of a value in the list and returns the index.

❷ `index` finds the *first* occurrence of a value in the list. In this case, `'new'` occurs twice in the list, in `li[2]` and `li[6]`, but `index` will return only the first index, `2`.

❸ If the value is not found in the list, Python raises an exception. This is notably different from most languages, which will return some invalid index. While this may seem annoying, it is a good thing, because it means your program will crash at the source of the problem, rather than later on when you try to use the invalid index.

❹ 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 (like an `if` statement), according to the following rules:

- 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; these values, like everything else in Python, are case–sensitive.

## 3.2.4. Deleting List Elements

**Example 3.13. Removing Elements from a List**

```
>>> li
['a', 'b', 'new', 'mpilgrim', 'z', 'example', 'new', 'two', 'elements']
>>> li.remove("z")    ❶
>>> li
['a', 'b', 'new', 'mpilgrim', 'example', 'new', 'two', 'elements']
>>> li.remove("new")  ❷
```

```
>>> li
['a', 'b', 'mpilgrim', 'example', 'new', 'two', 'elements']
>>> li.remove("c")    ❸
```

❹

❶
❷

❸
❹

❶

❷

❸

❶

❷

❸

```
>>> li
['a', 'b', 'mpilgrim', 'example', 'new', 'two', 'elements']
>>> li.remove("c")    ❸
```

- *Python Library Reference* (http://www.python.org/doc/current/lib/) summarizes all the list methods (http://www.python.org/doc/current/lib/typesseq−mutable.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**

```
>>> t = ("a", "b", "mpilgrim", "z", "example") ❶
>>> t
('a', 'b', 'mpilgrim', 'z', 'example')
>>> t[0]                                          ❷
'a'
>>> t[-1]                                         ❸
'example'
>>> t[1:3]                                        ❹
('b', 'mpilgrim')
```

❶    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.

❷    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 always `t[0]`.

❸    Negative indices count from the end of the tuple, just as with a list.

❹    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**

```
>>> t
('a', 'b', 'mpilgrim', 'z', 'example')
>>> t.append("new")      ❶
Traceback (innermost last):
  File "<interactive input>", line 1, in ?
AttributeError: 'tuple' object has no attribute 'append'
>>> t.remove("z")        ❷
Traceback (innermost last):
  File "<interactive input>", line 1, in ?
AttributeError: 'tuple' object has no attribute 'remove'
>>> t.index("example")   ❸
Traceback (innermost last):
  File "<interactive input>", line 1, in ?
AttributeError: 'tuple' object has no attribute 'index'
>>> "z" in t             ❹
True
```

❶    You can't add elements to a tuple. Tuples have no `append` or `extend` method.

❷    You can't remove elements from a tuple. Tuples have no `remove` or `pop` method.

❸    You can't find elements in a tuple. Tuples have no `index` method.

❹    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 instead of a list.

- 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 implied `assert` statement that shows this data is constant, and that special thought (and a specific function) is required to override that.
- 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 as keys in a dictionary, but lists can't be used this way.Actually, it's more complicated than that. Dictionary keys must be immutable. Tuples themselves are immutable, but if you have a tuple of lists, that counts as mutable and isn't safe to use as a dictionary key. Only tuples of strings, numbers, or other dictionary–safe tuples can be used as dictionary keys.
- Tuples are used in string formatting, as you'll see shortly.

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 `list` function takes a tuple and returns a list. In effect, `tuple` freezes a list, and `list` thaws a tuple.

**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 (http://www.ibiblio.org/obp/thinkCSpy/chap10.htm).
- Python Knowledge Base (http://www.faqts.com/knowledge–base/index.phtml/fid/199/) shows how to sort a tuple (http://www.faqts.com/knowledge–base/view.phtml/aid/4553/fid/587).
- *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/node7.html#SECTION007300000000000000000).

# 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, `odbchelper.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, and they are automatically destroyed when they go out of scope.

**Example 3.17. Defining the `myParams` Variable**

```
if __name__ == "__main__":
    myParams = {"server":"mpilgrim", \
                "database":"master", \
                "uid":"sa", \
                "pwd":"secret" \
                }
```

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's normally stringent indentation rules do not apply. If your Python IDE auto–indents the continued line, you should probably accept its default unless you have a burning reason not to.

Strictly speaking, expressions in parentheses, straight brackets, or curly braces (like defining a dictionary) can be split into multiple lines with or without the line continuation character ("\

❶

❶

❶
❷
❸

❶     The built–in `range` function returns a list of integers. In its simplest form, it takes an upper limit and returns a zero–based list counting up to but not including the upper limit. (If you like, you can pass other parameters to specify a base other than `0` and a step other than `1`. You can `print range.__doc__` for details.)

❷     `MONDAY`, `TUESDAY`, `WEDNESDAY`, `THURSDAY`, `FRIDAY`, `SATURDAY`, and `SUNDAY` are the variables you're defining. (This example came from the `calendar` module, a fun little module that prints calendars, like the UNIX program `cal`. The `calendar` module defines integer constants for days of the week.)

❸     Now each variable has its value: `MONDAY` is `0`, `TUESDAY` 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 caller can treat it as a tuple, or assign the values to individual variables. Many standard Python libraries do this, including the `os` module, which you'll discuss in Chapter 6.

### 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 (http://www.python.org/doc/current/ref/implicit–joining.html) and when you need to use it (http://www.python.org/doc/current/ref/explicit–joining.html).
- *How to Think Like a Computer Scientist* (http://www.ibiblio.org/obp/thinkCSpy/) shows how to use multi–variable assignment to swap the values of two variables (http://www.ibiblio.org/obp/thinkCSpy/chap09.htm).

## 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 a string with the `%s` placeholder.

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

### Example 3.21. Introducing String Formatting

```
>>> k = "uid"
>>> v = "sa"
>>> "%s=%s" % (k, v)   ❶
'uid=sa'
```

❶     The whole expression evaluates to a string. The first `%s` is replaced by the value of `k`; the second `%s` is replaced by the value of `v`. All other characters in the string (in this case, the equal sign) stay as they are.

Note that `(k, v)` is a tuple. I told you they were good for something.

You might be thinking that this is a lot of work just to do simple string concatentation, and you would be right, except that string formatting isn't just concatenation. It's not even just formatting. It's also type coercion.

### Example 3.22. String Formatting vs. Concatenating

```
>>> uid = "sa"
>>> pwd = "secret"
>>> print pwd + " is not a good password for " + uid        ❶
secret is not a good password for sa
>>> print "%s is not a good password for %s" % (pwd, uid)   ❷
```

```
secret is not a good password for sa
>>> userCount = 6
>>> print "Users connected: %d" % (userCount, )          ❸ ❹
Users connected: 6
>>> print "Users connected: " + userCount                ❺
Traceback (innermost last):
  File "<interactive input>", line 1, in ?
TypeError: cannot concatenate 'str' and 'int' objects
```

❶     + is the string concatenation operator.

❷     In this trivial case, string formatting accomplishes the same result as concatenation.

❸     `(userCount, )` is a tuple with one element. Yes, the syntax is a little strange, but there's a good reason for it: it's unambiguously a tuple. In fact, you can always include a comma after the last element when defining a list, tuple, or dictionary, but the comma is required when defining a tuple with one element. If the comma weren't required, Python wouldn't know whether `(userCount)` was a tuple with one element or just the value of `userCount`.

❹     String formatting works with integers by specifying `%d` instead of `%s`.

❺     Trying to concatenate a string with a non–string raises an exception. Unlike string formatting, string concatenation works only when everything is already a string.

As with `printf` in C, string formatting in Python is like a Swiss Army knife. There are options galore, and modifier strings to specially format many different types of values.

**Example 3.23. Formatting Numbers**

```
>>> print "Today's stock price: %f" % 50.4625    ❶
50.462500
>>> print "Today's stock price: %.2f" % 50.4625  ❷
50.46
>>> print "Change since yesterday: %+.2f" % 1.5  ❸
+1.50
```

❶     The `%f` string formatting option treats the value as a decimal, and prints it to six decimal places.

❷     The ".2" modifier of the `%f` option truncates the value to two decimal places.

❸     You can even combine modifiers. Adding the + modifier displays a plus or minus sign before the value. Note that the ".2" modifier is still in place, and is padding the value to exactly two decimal places.

**Further Reading on String Formatting**

- *Python Library Reference* (http://www.python.org/doc/current/lib/) summarizes all the string formatting format characters (http://www.python.org/doc/current/lib/typesseq–strings.html).
- *Effective AWK Programming* (http://www–gnats.gnu.org:8080/cgi–bin/info2www?(gawk)Top) discusses all the format characters (http://www–gnats.gnu.org:8080/cgi–bin/info2www?(gawk)Control+Letters) and advanced string formatting techniques like specifying width, precision, and zero–padding (http://www–gnats.gnu.org:8080/cgi–bin/info2www?(gawk)Format+Modifiers).

## 3.6. Mapping Lists

One of the most powerful features of Python is the list comprehenttp://6s81 –26.,chionesirt9l0ctiea–27ython ile stratp://chiones

```
>>> li = [1, 9, 8, 4]
>>> [elem*2 for elem in li]      ❶
[2, 18, 16, 8]
>>> li                           ❷
[1, 9, 8, 4]
>>> li = [elem*2 for elem in li] ❸
>>> li
[2, 18, 16, 8]
```

❶     To make sense of this, look at it from right to left. `li` is the list you're mapping. Python loops through `li` one element at a time, temporarily assigning the value of each element to the variable `elem`. Python then applies the function `elem*2` and appends that result to the returned list.

❷     Note that list comprehensions do not change the original list.

❸     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 when the list comprehension is complete, it assigns the result to the variable.

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

```
["%s=%s" % (k, v) for k, v in params.items()]
```

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

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

```
>>> params = {"server":"mpilgrim", "database":"master", "uid":"sa", "pwd":"secret"}
>>> params.keys()      ❶
['server', 'uid', 'database', 'pwd']
>>> params.values() ❷
['mpilgrim', 'sa', 'master', 'secret']
>>> params.items()     ❸
in the dictionarym'), ('uid', 'sa'), ('database', 'master'), ('pwd', 'secret')]27.65 Do Q > par . 
```

❶

**Exaare the list compreNoen et'ultee wed l0 −16 (2el 11 /F4 11 Tf ( function that you declardoedha retat 5. The )Tj /F5 11**

F5 Do Q18, 16, 8]'uid', 's
>>> params.values()

❸

❶

❷

```
['mpilgrim', 'sa', 'master', 'secret']
>>> ["%s=%s" % (k, v) for k, v in params.items()]  ❸
['server=mpilgrim', 'uid=sa', 'database=master', 'pwd=secret']
```

❶    Note that you're using two variables to iterate through the `params.items()` list. This is another use of multi–variable assignment. The first element of `params.items()` is `('server', 'mpilgrim')`, so in

❷

❸

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**

```
>>> li = ['server=mpilgrim', 'uid=sa', 'database=master', 'pwd=secret']
>>> s = ";".join(li)
>>> s
'server=mpilgrim;uid=sa;database=master;pwd=secret'
>>> s.split(";")          ❶
['server=mpilgrim', 'uid=sa', 'database=master', 'pwd=secret']
>>> s.split(";", 1) ❷
['server=mpilgrim', 'uid=sa;database=master;pwd=secret']
```

❶    `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 returned list.

❷    `split` takes an optional second argument, which is the number of times to split. (""Oooooh, optional
     arguments..." You'll learn how to do this in your own functions in the next chapter.)

*anystring*`.split(`*delimiter*`, 1)` 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 the first element of the returned list) and
everything after it (which ends up in the second element).

**Further Reading on String Methods**

- Python Knowledge Base (http://www.faqts.com/knowledge−base/index.phtml/fid/199/) answers common
  questions about strings (http://www.faqts.com/knowledge−base/index.phtml/fid/480) and has a lot of example
  code using strings (http://www.faqts.com/knowledge−base/index.phtml/fid/539).
- *Python Library Reference* (http://www.python.org/doc/current/lib/) summarizes all the string methods
  (http://www.python.org/doc/current/lib/string−methods.html).
- *Python Library Reference* (http://www.python.org/doc/current/lib/) documents the `string` module
  (http://www.python.org/doc/current/lib/module−string.html).
- *The Whole Python FAQ* (http://www.python.org/doc/FAQ.html) explains why `join` is a string method
  (http://www.python.org/cgi−bin/faqw.py?query=4.96&querytype=simple&casefold=yes&req=search) instead
  of a list method.

### 3.7.1. Historical Note on String Methods

When I first learned Python, I expected `join` to be a method of a list, which would take the delimiter as an argument.
Many people feel the same way, and there's a story behind the `join` method. Prior to Python 1.6, strings didn't have
all these useful methods. There was a separate `string` module that contained all the string functions; each function
took a string as its first argument. The functions were deemed important enough to put onto the strings themselves,
which made sense for functions like `lower`, `upper`, and `split`. But many hard−core Python programmers objected
to the new `join` method, arguing that it should be a method of the list instead, or that it shouldn't move at all but
simply stay a part of the old `string` module (which still has a lot of useful stuff in it). I use the new `join` method
exclusively, but you will see code written either way, and if it really bothers you, you can use the old `string.join`
function instead.

## 3.8. Summary

The `odbchelper.py` program and its output should now make perfect sense.

```
def buildConnectionString(params):
    """Build a connection string from a dictionary of parameters.
```

```
    Returns string."""
    return ";".join(["%s=%s" % (k, v) for k, v in params.items()])

if __name__ == "__main__":
    myParams = {"server":"mpilgrim", \
                "database":"master", \
                "uid":"sa", \
                "pwd":"secret" \
                }
    print buildConnectionString(myParams)
```

Here is the output of `odbchelper.py`:

```
server=mpilgrim;uid=sa;database=master;pwd=secret
```

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

- Using the Python IDE to test expressions interactively
- Writing Python programs and running them from within your IDE, or from the command line
- Importing modules and calling their functions
- Declaring functions and using `doc strings`, local variables, and proper indentation
- Defining dictionaries, tuples, and lists
- Accessing attributes and methods of any object, including strings, lists, dictionaries, functions, and modules
- Concatenating values through string formatting
- Mapping lists into other lists using list comprehensions
- Splitting strings into lists and joining lists into strings

❶ ❷

## 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`**

```
>>> str(1)              ❶
'1'
>>> horsemen = ['war', 'pestilence', 'famine']
>>> horsemen
['war', 'pestilence', 'famine']
>>> horsemen.append('Powerbuilder')
>>> str(horsemen)       ❷
"['war', 'pestilence', 'famine', 'Powerbuilder']"
>>> str(odbchelper)     ❸
"<module 'odbchelper' from 'c:\\docbook\\dip\\py\\odbchelper.py'>"
>>> str(None)           ❹
'None'
```

❶   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.

❷   However, `str` works on any object of any type. Here it works on a list which you've constructed in bits and pieces.

❸   `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.

❹   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 the `info` function, as you'll see shortly.

At the heart of the `info` function is the powerful `dir` function. `dir` returns a list of the attributes and methods of any object: modules, functions, strings, lists, dictionaries... pretty much anything.

**Example 4.7. Introducing `dir`**

```
>>> li = []
>>> dir(li)             ❶
['append', 'count', 'extend', 'index', 'insert',
'pop', 'remove', 'reverse', 'sort']
>>> d = {}
>>> dir(d)              ❷
['clear', 'copy', 'get', 'has_key', 'items', 'keys', 'setdefault', 'update', 'values']
>>> import odbchelper
>>> dir(odbchelper)     ❸
['__builtins__', '__doc__', '__file__', '__name__', 'buildConnectionString']
```

❶   `li` is a list, so `dir(li)` 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.

❷   `d` is a dictionary, so `dir(d)` returns a list of the names of dictionary methods. At least one of these, `keys`, should look familiar.

❸   This is where it really gets interesting. `odbchelper` is a module, so `dir(odbchelper)` returns a list of all kinds of stuff defined in the module, including built–in attributes, like `__name__`, `__doc__`, and whatever other attributes and methods you define. In this case, `odbchelper` has only one user–defined method, the `buildConnectionString` function described in Chapter 2.

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, even classes themselves. (More on classes in the next chapter.)

**Example 4.8. Introducing `callable`**

```
>>> import string
>>> string.punctuation            ❶
'!"#$%&\'()*+,-./:;<=>?@[\\]^_`{|}~'
>>> string.join                   ❷
<function join at 00C55A7C>
>>> callable(string.punctuation)  ❸
False
>>> callable(string.join)         ❹
True
>>> print string.join.__doc__     ❺
join(list [,sep]) -> string

    Return a string composed of the words in list, with
    intervening occurrences of sep.  The default separator is a
    single space.

    (joinfields and join are synonymous)
```

❶   The functions in the `string` module are deprecated (although many people still use the `join` function), but the module contains a lot of useful constants like this `string.punctuation`, which contains all the standard punctuation characters.

❷   `string.join` is a function that joins a list of strings.

❸   `string.punctuation` is not callable; it is a string. (A string does have callable methods, but the string itself is not callable.)

❹   `string.join` is callable; it's a function that takes two arguments.

❺   Any callable object may have a `doc string`. By using the `callable` function on each of an object's attributes, you can determine which attributes you care about (methods, functions, classes) and which you want to ignore (constants and so on) without knowing anything about the object ahead of time.

## 4.3.3. Built–In Functions

`type`, `str`, `dir`, and all the rest of Python's built–in functions are grouped into a special module called `__builtin__`. (That's two underscores before and after.) If it helps, you can think of Python automatically executing `from  __builtin__  import  *` on startup, which imports all the "built–in" functions into the namespace so you can use them directly.

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 `__builtin__` module. And guess what, Python has a function called `info`. Try it yourself and skim through the list now. We'll dive into some of the more important functions later. (Some of the built–in error classes, like `AttributeError`, should already look familiar.)

**Example 4.9. Built–in Attributes and Functions**

```
>>> from apihelper import info
>>> import __builtin__
>>> info(__builtin__, 20)
ArithmeticError      Base class for arithmetic errors.
AssertionError       Assertion failed.
AttributeError       Attribute not found.
```

EOFError          Read beyond end of file.

❶
❷
❸

❹
❺

❶

❷

❸

❹
❺

## 4.4.1. `getattr` with Modules

`getattr`

❶

❷

❸

❹

❺

❶

❷

❸

❹
❺

❶
❷
❸

❶ The `output` function takes one required argument, `data`, and one optional argument, `format`. If `format` is not specified, it defaults to `text`, and you will end up calling the plain text output function.

❷ You concatenate the `format` argument with "output_" to produce a function name, and then go get that function from the `statsout` module. This allows you to easily extend the program later to support other output formats, without changing this dispatch function. Just add another function to `statsout` named, for instance, `output_pdf`, and pass "pdf" as the `format` into the `output` function.

❸ Now you can simply call the output function in the same way as any other function. The `output_function` variable is a reference to the appropriate function from the `statsout` module.

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 the user passes in a format that doesn't have a corresponding function defined in `statsout`? Well, `getattr` will return `None`, which will be assigned to `output_function` instead of a valid function, and the next line that attempts to call that function will crash and raise an exception. That's bad.

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

**Example 4.13. `getattr` Default Values**

```
import statsout

def output(data, format="text"):
    output_function = getattr(statsout, "output_%s" % format, statsout.output_text)
    return output_function(data) ❶
```

❶ This function call is guaranteed to work, because you added a third argument to the call to `getattr`. The third argument is a default value that is returned if the attribute or method specified by the second argument wasn't found.

As you can see, `getattr` 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 can be combined with a filtering mechanism, where some elements in the list are mapped while others are skipped entirely.

Here is the list filtering syntax:

```
[mapping-expression for element in source-list if filter-expression]
```

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 filter expression. A filter expression can be any expression that evaluates true or false (which in Python can be almost anything). Any element for which the filter expression evaluates true will be included in the mapping. All other elements are ignored, so they are never put through the mapping expression and are not included in the output list.

**Example 4.14. Introducing List Filtering**

```
>>> li = ["a", "mpilgrim", "foo", "b", "c", "b", "d", "d"]
>>> [elem for elem in li if len(elem) > 1]        ❶
['mpilgrim', 'foo']
>>> [elem for elem in li if elem != "b"]          ❷
```

```
['a', 'mpilgrim', 'foo', 'c', 'd', 'd']
>>> [elem for elem in li if li.count(elem) == 1] ❸
['a', 'mpilgrim', 'foo', 'c']
```

❶     The mapping expression here is simple (it just returns the value of each element), so concentrate on the filter expression. As Python loops through the list, it runs each element through the filter expression. If the filter expression is true, the element is mapped and the result of the mapping expression is included in the returned list. Here, you are filtering out all the one–character strings, so you're left with a list of all the longer strings.

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

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 a list, returning a list containing only one copy of each value in the original list. But it doesn't, because values that appear twice in the original list (in this case, b and d) are excluded completely. There are ways of eliminating duplicates from a list, but filtering is not the solution.

Let's get back to this line from apihelper.py:

```
methodList = [method for method in dir(object) if callable(getattr(object, method))]
```

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 methodList variable. The first half of the expression is the list mapping part. The mapping expression is an identity expression, which it returns the value of each element. dir(object) returns a list of object's attributes and methods –– that's the list you're mapping. So the only new part is the filter expression after the if.

The filter expression looks scary, but it's not. You already know about callable, getattr, and in. As you saw in the previous section, the expression getattr(object, method) returns a function object if object is a module and method is the name of a function in that module.

So this expression takes an object (named object). Then it gets a list of the names of the object's attributes, methods, functions, and a few other things. Then it filters that list to weed out all the stuff that you don't care about. You do the weeding out by taking the name of each attribute/method/function and getting a reference to the real thing, via the getattr function. Then you check to see if that object is callable, which will be any methods and functions, both built–in (like the pop method of a list) and user–defined (like the buildConnectionString function of the odbchelper module). You don't care about other attributes, like the __name__ attribute that's built in to every module.

**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

❶

```
'b'
>>> '' and 'b'          ❷
''
>>>                     ❸
```

❶

❷

❸

❶
❷
❸

❹

❶

❷
❸

❹

❶

❷

❶

This syntax looks similar to the *bool* `? a : b` expression in C. The entire expression is evaluated from left to right, so the `and` is evaluated first. `1 and 'first'` evalutes to `'first'`, then `'first' or 'second'` evalutes to `'first'`.

❷ `0 and 'first'` evalutes to `False`, and then `0 or 'second'` evaluates to `'second'`.

However, since this Python expression is simply boolean logic, and not a special construct of the language, there is one extremely important difference between this `and-or` trick in Python and the *bool* `? a : b` syntax in C. If the value of a is false, the expression will not work as you would expect it to. (Can you tell I was bitten by this? More than once?)

**Example 4.18. When the `and-or` Trick Fails**

```
>>> a = ""
>>> b = "second"
>>> 1 and a or b          ❶
'second'
```

❶ Since a is an empty string, which Python considers false in a boolean context, `1 and ''` evalutes to `''`, and then `'' or 'second'` evalutes to `'second'`. Oops! That's not what you wanted.

The `and-or` trick, *bool* `and a or b`, will not work like the C expression *bool* `? a : b` when a is false in a boolean context.

The real trick behind the `and-or` 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 `[a]` and b into `[b]`, then taking the first element of the returned list, which will be either a or b.

**Example 4.19. Using the `and-or` Trick Safely**

```
>>> a = ""
>>> b = "second"
>>> (1 and [a] or [b])[0] ❶
''
```

❶ Since `[a]` is a non−empty list, it is never false. Even if a is `0` or `''` or some other false value, the list `[a]` 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 fuss? Well, in many cases, you are choosing between two constant values, so you can use the simpler syntax and not worry, because you know that the a value will always be true. And even if you need to use the more complicated safe form, there are good reasons to do so. For example, there are some cases in Python where `if` statements are not allowed, such as in `lambda` functions.

**Further Reading on 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/52310).

## 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 anywhere a function is required.

## Example 4.20. Introducing `lambda` Functions

```
>>> def f(x):
...     return x*2
...
>>> f(3)
6
>>> g = lambda x: x*2      ❶
>>> g(3)
6
>>> (lambda x: x*2)(3)  ❷
6
```

❶    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 argument list, and the `return` keyword is missing (it is implied, since the entire function can only be one expression). Also, the function has no name, but it can be called through the variable it is assigned to.

❷    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 that a lambda is just an in–line function.

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 can not contain commands, and they can not contain more than one expression. Don't try to squeeze too much into a `lambda` function; if you need something more complex, define a normal function instead and make it as long as you want.

`lambda` functions are a matter of style. Using them is nariection canlvul/F4 nFer owno pa

❶

❷

❸

❶

❶
❷
❸

❹

❶
❷
❸
❹

❶
❷

❸

❶

❷
❸

Stepping back even further, you see that you're using string formatting again to concatenate the return value of `processFunc` with the return value of `method`'s `ljust` method. This is a new string method that you haven't seen before.

**Example 4.24. Introducing `ljust`**

```
>>> s = 'buildConnectionString'
>>> s.ljust(30) ❶
'buildConnectionString         '
>>> s.ljust(20) ❷
'buildConnectionString'
```

❶     `ljust` pads the string with spaces to the given length. This is what the `info` function uses to make two columns of output and line up all the `doc strings` in the second column.

❷     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 `ljust` method and the (possibly collapsed) `doc string` from the call to `processFunc`

llapsed ged. It

>>> s = 'buildConn5. PssFn4 11  u will simply return the string unchang',is w−16.276 Tmk4 11  u will simply return the string uncha
columns of /F4 11 Tfs of TConn4.50 0.50 e /mrinf (doc string)Tj f ( will sg.)Tj −2just

❶

❶

```
      print info.__doc__
```

Here is the output of `apihelper.py`:

```
>>> from apihelper import info
>>> li = []
>>> info(li)
append      L.append(object) -- append object to end
count       L.count(value) -> integer -- return number of occurrences of value
extend      L.extend(list) -- extend list by appending list elements
index       L.index(value) -> integer -- return index of first occurrence of value
insert      L.insert(index, object) -- insert object before index
pop         L.pop([index]) -> item -- remove and return item at index (default last)
remove      L.remove(value) -- remove first occurrence of value
reverse     L.reverse() -- reverse *IN PLACE*
sort        L.sort([cmpfunc]) -- sort *IN PLACE*; if given, cmpfunc(x, y) -> -1, 0, 1
```

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 `and-or` trick and using it safely
- Defining `lambda` functions
- Assigning functions to variables and calling the function by referencing the variable. I can't emphasize this enough, because this mode of thought is vital to advancing your understanding of Python. You'll see more complex applications of this concept throughout this book.

# 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 strings` of the module, the classes, and the functions to get an overview of what this program does and how it works. As usual, don't worry about the stuff you don't understand; that's what the rest of the chapter is for.

**Example 5.1. `fileinfo.py`**

If you have not already done so, you can download this and other examples
(http://diveintopython.org/download/diveintopython–examples–5.4.zip) used in this book.

```
"""Framework for getting filetype-specific metadata.

Instantiate appropriate class with filename.  Returned object acts like a
dictionary, with key-value pairs for each piece of metadata.
    import fileinfo
    info = fileinfo.MP3FileInfo("/music/ap/mahadeva.mp3")
    print "\\n".join(["%s=%s" % (k, v) for k, v in info.items()])

Or use listDirectory function to get info on all files in a directory.
    for info in fileinfo.listDirectory("/music/ap/", [".mp3"]):
        ...

Framework can be extended by adding classes for particular file types, e.g.
HTMLFileInfo, MPGFileInfo, DOCFileInfo.  Each class is completely responsible for
parsing its files appropriately; see MP3FileInfo for example.
"""
import os
import sys
from UserDict import UserDict

def stripnulls(data):
    "strip whitespace and nulls"
    return data.replace("\00", "").strip()

class FileInfo(UserDict):
    "store file metadata"
    def __init__(self, filename=None):
        UserDict.__init__(self)
        self["name"] = filename

class MP3FileInfo(FileInfo):
    "store ID3v1.0 MP3 tags"
    tagDataMap = {"title"   : (  3,  33, stripnulls),
                  "artist"  : ( 33,  63, stripnulls),
                  "album"   : ( 63,  93, stripnulls),
                  "year"    : ( 93,  97, stripnulls),
                  "comment" : ( 97, 126, stripnulls),
                  "genre"   : (127, 128, ord)}

    def __parse(self, filename):
        "parse ID3v1.0 tags from MP3 file"
        self.clear()
        try:
```

```
                fsock = open(filename, "rb", 0)
                try:
                    fsock.seek(-128, 2)
                    tagdata = fsock.read(128)
                finally:
                    fsock.close()
                if tagdata[:3] == "TAG":
                    for tag, (start, end, parseFunc) in self.tagDataMap.items():
                        self[tag] = parseFunc(tagdata[start:end])
            except IOError:
                pass

        def __setitem__(self, key, item):
            if key == "name" and item:
                self.__parse(item)
            FileInfo.__setitem__(self, key, item)

def listDirectory(directory, fileExtList):
    "get list of file info objects for files of particular extensions"
    fileList = [os.path.normcase(f)
                for f in os.listdir(directory)]
    fileList = [os.path.join(directory, f)
                for f in fileList
                if os.path.splitext(f)[1] in fileExtList]
    def getFileInfoClass(filename, module=sys.modules[FileInfo.__module__]):
        "get file info class from filename extension"
        subclass = "%sFileInfo" % os.path.splitext(filename)[1].upper()[1:]
        return hasattr(module, subclass) and getattr(module, subclass) or FileInfo
    return [getFileInfoClass(f)(f) for f in fileList]

if __name__ == "__main__":
    for info in listDirectory("/music/_singles/", [".mp3"]):  ❶
        print "\n".join(["%s=%s" % (k, v) for k, v in info.items()])
        print
```

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 point to a directory of MP3 files on your own machine.

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 music.

```
album=
artist=Ghost in the Machine
title=A Time Long Forgotten (Concept
genre=31
name=/music/_singles/a_time_long_forgotten_con.mp3
year=1999
comment=http://mp3.com/ghostmachine

album=Rave Mix
artist=***DJ MARY-JANE***
title=HELLRAISER****Trance from Hell
genre=31
name=/music/_singles/hellraiser.mp3
year=2000
comment=http://mp3.com/DJMARYJANE

album=Rave Mix
artist=***DJ MARY-JANE***
title=KAIRO****THE BEST GOA
genre=31
name=/music/_singles/kairo.mp3
year=2000
```

```
comment=http://mp3.com/DJMARYJANE

album=Journeys
artist=Masters of Balance
title=Long Way Home
genre=31
name=/music/_singles/long_way_home1.mp3
year=2000
comment=http://mp3.com/MastersofBalan

album=
artist=The Cynic Project
title=Sidewinder
genre=18
name=/music/_singles/sidewinder.mp3
year=2000
comment=http://mp3.com/cynicproject

album=Digitosis@128k
artist=VXpanded
title=Spinning
genre=255
name=/music/_singles/spinning.mp3
year=2000
comment=http://mp3.com/artists/95/vxp
```

## 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, `import module`, you've already seen in Section 2.4,  Everything Is an Object . The other way accomplishes the same thing, but it has subtle and important differences.

Here is the basic `from module import` syntax:

```
from UserDict import UserDict
```

This is similar to the `import module` syntax that you know and love, but with an important difference: the attributes and methods of the imported module `types` are imported directly into the local namespace, so they are available directly, without qualification by module name. You can import individual items or use `from module import *` to import everything.

`from module import *` in Python is like `use module` in Perl; `import module` in Python is like `require module` in Perl.

`from module import *` in Python is like `import module.*` in Java; `import module` in Python is like `import module` in Java.

**Example 5.2. `import module` *vs.* `from module import`**

```
>>> import types
>>> types.FunctionType                    ❶

                                          ❷



                                          ❸
```

```
>>> FunctionType                    ❹
<type 'function'>
```

❶     The `types` module contains no methods; it just has attributes for each Python object type. Note that the attribute, `FunctionType`, must be qualified by the module name, `types`.

❷     `FunctionType` by itself has not been defined in this namespace; it exists only in the context of `types`.

❸     This syntax imports the attribute `FunctionType` from the `types` module directly into the local namespace.

❹     Now `FunctionType` can be accessed directly, without reference to `types`.

When should you use `from module import`?

- If you will be accessing attributes and methods often and don't want to type the module name over and over, use `from module import`.
- If you want to selectively import some attributes and methods but not others, use `from module import`.
- If the module contains attributes or functions with the same name as ones in your module, you must use `import module` 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 module import` * sparingly, because it makes it difficult to determine where a particular function or attribute came from, and that makes debugging and refactoring more difficult.

**Further Reading on Module Importing Techniques**

- eff–bot (http://www.effbot.org/guides/) has more to say on `import module` *vs.* `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/node8.html#SECTION008410000000000000000).

# 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.

       ❶
      ❷ ❸

❶

❷

**❸**      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 first thing not indented is not in the class.

The `pass` 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've just seen, there is nothing that a class absolutely must have, other than a name. In particular, C++ programmers may find it odd that Python classes don't have explicit constructors and destructors. Python classes do have something similar to a constructor: the __init__ method.

**Example 5.4. Defining the `FileInfo` Class**

```
from UserDict import UserDict

class FileInfo(UserDict):   ❶
```

**❶**      In Python, the ancestor of a class is simply listed in parentheses immediately after the class name. So the `FileInfo` class is inherited from the `UserDict` class (which was imported from the `UserDict` module). `UserDict` is a class that acts like a dictionary, allowing you to essentially subclass the dictionary datatype and add your own behavior. (There are similar classes `UserList` and `UserString` which allow you to subclass lists and strings.) There is a bit of black magic behind this, which you will demystify later in this chapter when you explore the `UserDict` class in more depth.

In Python, the ancestor of a class is simply listed in parentheses immediately after the class name. There is no special keyword like `extends` 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 `FileInfo` class using the __init__ method.

**Example 5.5. Initializing the `FileInfo` Class**

<div align="center">

❶
❷   ❸   ❹

</div>

**❶**
**❷**




**❸**

## 5.4. Instantiating Classes

❶
❷

❸

❹

❶

❷

❸

❹

☞

❶

❷

❶

**❸**   Python supports data attributes (called "instance variables" in Java and Powerbuilder, and "member variables" in C++). Data attributes are pieces of data held by a specific instance of a class. In this case, each instance of `UserDict` will have a data attribute `data`. To reference this attribute from code outside the class, you qualify it with the instance name, *instance*`.data`, in the same way that you qualify a function with its module name. To reference a data attribute from within the class, you use `self` as the qualifier. By convention, all data attributes are initialized to reasonable values in the `__init__` method. However, this is not required, since data attributes, like local variables, spring into existence when they are first assigned a value.

**❹**   The `update` method is a dictionary duplicator: it copies all the keys and values from one dictionary to another. This does *not* clear the target dictionary first; if the target dictionary already has some keys, the ones from the source dictionary will be overwritten, but others will be left untouched. Think of `update` as a merge function, not a copy function.

**❺**   This is a syntax you may not have seen before (I haven't used it in the examples in this book). It's an `if` statement, but instead of having an indented block starting on the next line, there is just a single statement on the same line, after the colon. This is perfectly legal syntax, which is just a shortcut you can use when you have only one statement in a block. (It's like specifying a single statement without braces in C++.) You can use this syntax, or you can have indented code on subsequent lines, but you can't do both for the same block.

Java and Powerbuilder support function overloading by argument list, *i.e.* one class can have multiple methods with the same name but a different number of arguments, or arguments of different types. Other languages (most notably PL/SQL) even support function overloading by argument name; *i.e.* one class can have multiple methods with the same name and the same number of arguments of the same type but different argument names. Python supports neither of these; it has no form of function overloading whatsoever. Methods are defined solely by their name, and there can be only one method per class with a given name. So if a descendant class has an `__init__` method, it *always* overrides the ancestor `__init__` method, even if the descendant defines it with a different argument list. And the same rule applies to any other method.

Guido, the original author of Python, explains method overriding this way: "Derived classes may override methods of their base classes. Because methods have no special privileges when calling other methods of the same object, a method of a base class that calls another method defined in the same base class, may in fact end up calling a method of a derived class that overrides it. (For C++ programmers: all methods in Python are effectively virtual.)" If that doesn't make sense to you (it confuses the hell out of me), feel free to ignore it. I just thought I'd pass it along.

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 `AttributeError` exceptions because you're referencing uninitialized (and therefore non−existent) attributes.

### Example 5.10. `UserDict` Normal Methods

```
def clear(self): self.data.clear()              ❶
def copy(self):                                 ❷
    if self.__class__ is UserDict:              ❸
        return UserDict(self.data)
    import copy                                 ❹
    return copy.copy(self)
def keys(self): return self.data.keys()         ❺
def items(self): return self.data.items()
def values(self): return self.data.values()
```

**❶**   `clear`

❷ The `copy` method of a real dictionary returns a new dictionary that is an exact duplicate of the original (all the same key–value pairs). But `UserDict` can't simply redirect to `self.data.copy`, because that method returns a real dictionary, and what you want is to return a new instance that is the same class as `self`.

❸ You use the `__class__` attribute to see if `self` is a `UserDict`; if so, you're golden, because you know how to copy a `UserDict`: just create a new `UserDict` and give it the real dictionary that you've squirreled away in `self.data`. Then you immediately return the new `UserDict` you don't even get to the `import copy` on the next line.

❹ If `self.__class__` is not `UserDict`, then `self` must be some subclass of `UserDict` (like maybe `FileInfo`), in which case life gets trickier. `UserDict` doesn't know how to make an exact copy of one of its descendants; there could, for instance, be other data attributes defined in the subclass, so you would need to iterate through them and make sure to copy all of them. Luckily, Python comes with a module to do exactly this, and it's called `copy`. I won't go into the details here (though it's a wicked cool module, if you're ever inclined to dive into it on your own). Suffice it to say that `copy` can copy arbitrary Python objects, and that's how you're using it here.

❺ The rest of the methods are straightforward, redirecting the calls to the built–in methods on `self.data`.

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 comes with wrapper classes that mimic the behavior of these built–in datatypes: `UserString`, `UserList`, and `UserDict`. Using a combination of normal and special methods, the `UserDict` class does an excellent imitation of a dictionary. In Python 2.2 and later, you can inherit classes directly from built–in datatypes like `dict`. An example of this is given in the examples that come with this book, in `fileinfo_fromdict.py`.

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

**Example 5.11. Inheriting Directly from Built–In Datatype `dict`**

```
class FileInfo(dict):                    ❶
    "store file metadata"
    def __init__(self, filename=None):   ❷
        self["name"] = filename
```

❶ The first difference is that you don't need to import the `UserDict` module, since `dict` is a built–in datatype and is always available. The second is that you are inheriting from `dict` directly, instead of from `UserDict.UserDict`.

❷ The third difference is subtle but important. Because of the way `UserDict` works internally, it requires you to manually call its `__init__` method to properly initialize its internal data structures. `dict` does not work like this; it is not a wrapper, and it requires no explicit initialization.

**Further Reading on `UserDict`**

- *Python Library Reference* (http://www.python.org/doc/current/lib/) documents the `UserDict` module (http://www.python.org/doc/current/lib/module–UserDict.html) and the `copy` module (http://www.python.org/doc/current/lib/module–copy.html).

# 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 your code (like normal methods), special methods are called for you by Python in particular circumstances or when specific syntax is used.

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, because there are a lot of things you can do with dictionaries besides call methods on them. For starters, you can get and set items with a syntax that doesn't include explicitly invoking methods. This is where special class methods come in: they provide a way to map non–method–calling syntax into method calls.

## 5.6.1. Getting and Setting Items

**Example 5.12. The __getitem__ Special Method**

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

>>> f = fileinfo.FileInfo("/music/_singles/kairo.mp3")
>>> f
{'name':'/music/_singles/kairo.mp3'}
>>> f.__getitem__("name")  ❶
'/music/_singles/kairo.mp3'
>>> f["name"]              ❷
'/music/_singles/kairo.mp3'
```

❶    The __getitem__ special method looks simple enough. Like the normal methods '/muo_singles/kairo.mp3'

❷

❶

❷

❶

❷

___setitem___ is a special class method because it gets called for you, but it's still a class method. Just as easily as the ___setitem___ method was defined in UserDict, you can redefine it in the descendant class to override the ancestor method. This allows you to define classes that act like dictionaries in some ways but define their own behavior above and beyond the built−in dictionary.

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 metadata from a particular type of file. Once some attributes (like the file's name and location) are known, the handler class knows how to derive other attributes automatically. This is done by overriding the ___setitem___ method, checking for particular keys, and adding additional processing when they are found.

For example, MP3FileInfo is a descendant of FileInfo. When an MP3FileInfo's name is set, it doesn't just set the name key (like the ancestor FileInfo does); it also looks in the file itself for MP3 tags and populates a whole set of keys. The next example shows how this works.

**Example 5.14. Overriding ___setitem___ in MP3FileInfo**

```
def __setitem__(self, key, item):          ❶
    if key == "name" and item:             ❷
        self.__parse(item)                 ❸
    FileInfo.__setitem__(self, key, item)  ❹
```

❶  Notice that this ___setitem___ method is defined exactly the same way as the ancestor method. This is important, since Python will be calling the method for you, and it expects it to be defined with a certain number of arguments. (Technically speaking, the names of the arguments don't matter; only the number of arguments is important.)

❷  Here's the cru−13.2 Td(important.))Tj 0 M6 (__setitem__)T in ow to : −13.2 Td(imporis omatics easily as)itemnts. (Tech

❸

❹

❶

❷

❸

❶

❷

❸

❶
❷

❸
❹

❶

❷

❸

and code the length calculation yourself, and then call `len(`*`instance`*`)` and Python will call your `__len__` special method for you.

❹     `__delitem__` is called when you call `del` *`instance`*`[`*`key`*`]`, which you may remember as the way to delete individual items from a dictionary. When you use `del` on a class instance, Python calls the `__delitem__` special method for you.

In Java, you determine whether two string variables reference the same physical memory location by using `str1 == str2`. This is called *object identity*, and it is written in Python as `str1 is str2`. To compare string values in Java, you would use `str1.equals(str2)`; in Python, you would use `str1 == str2`. Java programmers who have been taught to believe that the world is a better place because `==` in Java compares by identity instead of by value may have a difficult time adjusting to Python's lack of such "gotchas".

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 would be easier (and the entire `UserDict` class would be unnecessary) if you could inherit from built−in datatypes like a dictionary. But even if you could, special methods would still be useful, because they can be used in any class, not just wrapper classes like `UserDict`.

Special methods mean that *any class* can store key/value pairs like a dictionary, just by defining the `__setitem__` method. *Any class* can act like a sequence, just by defining the `__getitem__` method. Any class that defines the `__cmp__` method can be compared with `==`. And if your class represents something that has a length, don't define a `GetLength` method; define the `__len__` method and use `len(`*`instance`*`)`.

While other object−oriented languages only let you define the physical model of an object ("this object has a `GetLength` method"), Python's special class methods like `__len__` allow you to define the logical model of an object ("this object has a length").

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 other arithmetic operations on class instances. (The canonical example of this is a class that represents complex numbers, numbers with both real and imaginary components.) The `__call__` method lets a class act like a function, allowing you to call a class instance directly. And there are other special methods that allow classes to have read−only and write−only data attributes; you'll talk more about those in later chapters.

**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/doc/current/ref/specialnames.html).

# 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 are variables owned by the class itself.

**Example 5.17. Introducing Class Attributes**

```
class MP3FileInfo(FileInfo):
    "store ID3v1.0 MP3 tags"
    tagDataMap = {"title"   : (  3,  33, stripnulls),
                  "artist"  : ( 33,  63, stripnulls),
                  "album"   : ( 63,  93, stripnulls),
                  "year"    : ( 93,  97, stripnulls),
                  "comment" : ( 97, 126, stripnulls),
                  "genre"   : (127, 128, ord)}
```

```
>>> import fileinfo
>>> fileinfo.MP3FileInfo                    ❶
<class fileinfo.MP3FileInfo at 01257FDC>
>>> fileinfo.MP3FileInfo.tagDataMap         ❷
{'title': (3, 33, <function stripnulls at 0260C8D4>),
 'genre': (127, 128, <built-in function ord>),
 'artist': (33, 63, <function stripnulls at 0260C8D4>),
 'year': (93, 97, <function stripnulls at 0260C8D4>),
 'comment': (97, 126, <function stripnulls at 0260C8D4>),
 'album': (63, 93, <function stripnulls at 0260C8D4>)}
>>> m = fileinfo.MP3FileInfo()              ❸
>>> m.tagDataMap
{'title': (3, 33, <function stripnulls at 0260C8D4>),
 'genre': (127, 128, <built-in function ord>),
 'artist': (33, 63, <function stripnulls at 0260C8D4>),
 'year': (93, 97, <function stripnulls at 0260C8D4>),
 'comment': (97, 126, <function stripnulls at 0260C8D4>),
 'album': (63, 93, <function stripnulls at 0260C8D4>)}
```

❶    MP3FileInfo is the class itself, not any particular instance of the class.

❷    tagDataMap is a class attribute: literally, an attribute of the class. It is available before creating any instances of the class.

❸    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 after the class definition (one with the static keyword, one without). In Python, only class attributes can be defined here; data attributes are defined in the __init__ method.

Class attributes can be used as class–level constants (which is how you use them in MP3FileInfo), 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 behavior should be discouraged but not banned. If you really want to change the value of None, you can do it, but don't come running to me when your code is impossible to debug.

**Example 5.18. Modifying Class Attributes**

```
>>> class counter:
...      count = 0                           ❶
...      def __init__(self):
...          self.__class__.count += 1       ❷
...
>>> counter
<class __main__.counter at 010EAECC>
>>> counter.count                            ❸
0
>>> c = counter()
>>> c.count                                  ❹
1
>>> counter.count
1
>>> d = counter()                            ❺
>>> d.count
2
>>> c.count
2
>>> counter.count
```

❶   count is a class attribute of the counter class.

❷   __class__ is a built–in attribute of every class instance (of every class). It is a reference to the class that self is an instance of (in this case, the counter class).

❸   Because count is a class attribute, it is available through direct reference to the class, before you have created any instances of the class.

❹   Creating an instance of the class calls the __init__ method, which increments the class attribute count by 1. This affects the class itself, not just the newly created instance.

❺   Creating a second instance will increment the class attribute count again. Notice how the class attribute is shared by the class and all instances of the class.

## 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. Python has no concept of *protected* class methods (accessible only in their own class and descendant classes). Class methods are either private (accessible only in their own class) or public (accessible from anywhere).

In MP3FileInfo, there are two methods: __parse and __setitem__. As you have already discussed, __setitem__ is a special method; normally, you would call it indirectly by using the dictionary syntax on a class instance, but it is public, and you could call it directly (even from outside the fileinfo module) if you had a really good reason. However, __parse is private, because it has two underscores at the beginning of its name.

In Python, all special methods (like __setitem__) and built–in attributes (like __doc__) follow a standard naming convention: they both start with and end with two underscores. Don't name your own methods and attributes w1 Tf t0eas o1.06nterf t0eass am. remebuht.tsho i –4mset>> 11 T6iom outMt –y, you wouca /_ ( mles/kairo.mp3".tsho mse4 1

❶

❶

the name `_MP3FileInfo__parse`

for an `IOError` exception, Python just prints out some debugging information about what happened and then gives up.

❷　You're trying to open the same non−existent file, but this time you're doing it within a `try...except` block.

❸　When the `open` method raises an `IOError` exception, you're ready for it. The `except IOError:` line catches the exception and executes your own block of code, which in this case just prints a more pleasant error message.

❹　Once an exception has been handled, processing continues normally on the first line after the `try...except` block. Note that this line will always print, whether or not an exception occurs. If you really did have a file called `notthere` in your root directory, the call to `open` would succeed, the `except` clause would be ignored, and this line would still be executed.

Exceptions may seem unfriendly (after all, if you don't catch the exception, your entire program will crash), but consider the alternative. Would you rather get back an unusable file object to a non−existent file? You'd need to check its validity somehow anyway, and if you forgot, somewhere down the line, your program would give you strange errors somewhere down the line that you would need to trace back to the source. I'm sure you've experienced this, and

❶

❷

❸

❹
❺

❶     `termios` is a UNIX–specific module that provides low–level control over the input terminal. If this module is not available (because it's not on your system, or your system doesn't support it), the import fails and Python raises an `ImportError`, which you catch.

❷     OK, you didn't have `termios`

❸

❹

❺

❶
❷

❸

❹

❶

❷

❸
❹

**Example 6.5. Closing a File**

```
>>> f
<open file '/music/_singles/kairo.mp3', mode 'rb' at 010E3988>
>>> f.closed            ❶
False
>>> f.close()           ❷
>>> f
<closed file '/music/_singles/kairo.mp3', mode 'rb' at 010E3988>
>>> f.closed            ❸
True
>>> f.seek(0)           ❹
Traceback (innermost last):
  File "<interactive input>", line 1, in ?
ValueError: I/O operation on closed file
>>> f.tell()
Traceback (innermost last):
  File "<interactive input>", line 1, in ?
ValueError: I/O operation on closed file
>>> f.read()
Traceback (innermost last):
  File "<interactive input>", line 1, in ?
ValueError: I/O operation on closed file
>>> f.close()           ❺
```

❶    The `closed`

❷

❸
❹

❺

❶
❷

❸
❹
❺

❻

❶ Because opening and reading files is risky and may raise an exception, all of this code is wrapped in a `try...except` block. (Hey, isn't standardized indentation great? This is where you start to appreciate it.)

❷ The `open` function may raise an `IOError`. (Maybe the file doesn't exist.)

❸ The `seek` method may raise an `IOError`. (Maybe the file is smaller than 128 bytes.)

❹ The `read` method may raise an `IOError`. (Maybe the disk has a bad sector, or it's on a network drive and the network just went down.)

❺ This is new: a `try...finally` block. Once the file has been opened successfully by the `open` function, you want to make absolutely sure that you close it, even if an exception is raised by the `seek` or `read` methods. That's what a `try...finally` block is for: code in the `finally` block will *always* be executed, even if something in the `try` block raises an exception. Think of it as code that gets executed on the way out, regardless of what happened before.

❻ At last, you handle your `IOError` exception. This could be the `IOError` exception raised by the call to `open`, `seek`, or `read`. Here, you really don't care, because all you're going to do is ignore it silently and continue. (Remember, `pass` is a Python statement that does nothing.) That's perfectly legal; "handling" an exception can mean explicitly doing nothing. It still counts as handled, and processing will continue normally on the next line of code after the `try...except` block.

## 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.
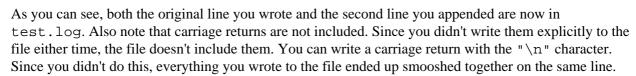- "write" mode will overwrite the file.

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 yet, create a new empty file just so you can open it for the first time" logic. Just open it and start writing.

**Example 6.7. Writing to Files**

```
>>> logfile = open('test.log', 'w') ❶
>>> logfile.write('test succeeded') ❷
>>> logfile.close()
>>> print file('test.log').read()    ❸
test succeeded
>>> logfile = open('test.log', 'a') ❹
>>> logfile.write('line 2')
>>> logfile.close()
>>> print file('test.log').read()    ❺
test succeededline 2
```

❶ You start boldly by creating either the new file `test.log` or overwrites the existing file, and opening the file for writing. (The second parameter `"w"` means open the file for writing.) Yes, that's all as dangerous as it sounds. I hope you didn't care about the previous contents of that file, because it's gone now.

❷ You can add data to the newly opened file with the `write` method of the file object returned by `open`.

❸ `file` is a synonym for `open`. This one–liner opens the file, reads its contents, and prints them.

❹ You happen to know that `test.log` exists (since you just finished writing to it), so you can open it and append to it. (The `"a"` parameter means open the file for appending.) Actually you could do this even if the file didn't exist, because opening the file for appending will create the file if necessary. But appending

will *never* harm the existing contents of the file.

❺  As you can see, both the original line you wrote and the second line you appended are now in
  `test.log`. Also note that carriage returns are not included. Since you didn't write them explicitly to the
  file either time, the file doesn't include them. You can write a carriage return with the `"\n"` character.
  Since you didn't do this, everything you wrote to the file ended up smooshed together on the same line.

**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 one line at a time into a list
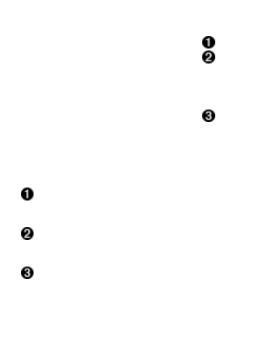  (http://www.python.org/doc/current/tut/node9.html#SECTION009210000000000000000).
- eff−bot (http://www.effbot.org/guides/) discusses efficiency and performance of various ways of reading a file
  (http://www.effbot.org/guides/readline−performance.htm).
- Python Knowledge Base (http://www.faqts.com/knowledge−base/index.phtml/fid/199/) answers common
  questions about files (http://www.faqts.com/knowledge−base/index.phtml/fid/552).
- *Python Library Reference* (http://www.python.org/doc/current/lib/) summarizes all the file object methods
  (http://www.python.org/doc/current/lib/bltin−file−objects.html).

# 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 many other things that you don't need them as often.

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

❶
❷


❸



❶

❷

❸



❶

```
0
1
2
3
4
>>> li = ['a', 'b', 'c', 'd', 'e']
>>> for i in range(len(li)):       ❷
...     print li[i]
a
b
c
d
e
```

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 a bit odd, but it is occasionally (and I stress *occasionally*) useful to have a
counter loop.

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.

`for` 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**

```
>>> import os
>>> for k, v in os.environ.items():       ❶ ❷
...     print "%s=%s" % thc',l)
USERPROFILE=C:\Documents and Settings\mpilgrim
OS=Windows_NT
COMPUTERNAME=MPILGRIM
USERNAME=mpilgrim

[...snip...]
>>> print "\n".join(["%s=%s" % thc',l)
...     for k, v in os.environ.items()])  ❸
USERPROFILE=C:\Documents and Settings\mpilgrim
OS=Windows_NT
COMPUTERNAME=MPILGRIM
USERNAME=mpilgrim

[...snip...]
```

❶    `os.environ` is a dictionary of the environment variables defined on your system. In Windows, these are your
user and system variables accessible from MS−DOS. In UNIX, they are the variables exported in your shell's
startup scripts. In Mac OS, there is no concept of environment variables, so this dictionary is empty.

❷    `os.environ.items()` returns a list of tuples: `[(key1, value1), (key2, value2), ...]`. The
`for` loop iterates through this list. The first round, it assigns *key1* to k and *value1* to v, so k =
`USERPROFILE` and v = `C:\Documents and Settings\mpilgrim`. In the second round, k gets the
second key, `OS`, and v gets the corresponding value, `Windows_NT`.

❸    With multi−variable assignment and list comprehensions, you can replace the entire `for` loop with a single
statement. Whether you actually do this in real code is a matter of personal coding style. I like it because it
makes it clear that what I'm doing is mapping a dictionary into a list, then joining the list into a single string.
Other programmers prefer to write this out as a `for` loop. The output is the same in either case, although this
version is slightly faster, because there is only one `print` statement instead of many.

Now we can look at the `for` loop in `MP3FileInfo`, from the sample `fileinfo.py` program introduced in

**Example 6.11. `for` Loop in `MP3FileInfo`**

```
tagDataMap = {"title"   : (  3,  33, stripnulls),
              "artist"  : ( 33,  63, stripnulls),
              "album"   : ( 63,  93, stripnulls),
              "year"    : ( 93,  97, stripnulls),
              "comment" : ( 97, 126, stripnulls),
              "genre"   : (127, 128, ord)}                    ❶
     .
     .
     .
        if tagdata[:3] == "TAG":
            for tag, (start, end, parseFunc) in self.tagDataMap.items():  ❷
                self[tag] = parseFunc(tagdata[start:end])                 ❸
```

❶    `tagDataMap` is 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 128 bytes of the file, bytes 3 through 32 of those are always the song title, 33 through 62 are always the artist name, 63 through 92 are the album name, and so forth. Note that `tagDataMap` is a dictionary of tuples, and each tuple contains two integers and a function reference.

❷    This looks complicated, but it's not. The structure of the `for` variables matches the structure of the elements of the list returned by `items`. Remember that `items` returns a list of tuples of the form (*key*, *value*). The

❸

 

 

❶
❷

 

❶

running (`sys.version` or `sys.version_info`), and system–level options such as the maximum allowed recursion depth (`sys.getrecursionlimit()` and `sys.setrecursionlimit()`).

❷  `sys.modules` is a dictionary containing all the modules that have ever been imported since

❶

❷

❶

❷

❶
❷

❶

❷

❶

❷
❸

❶

❷

❸

❶ ❷

❸

❹

❺

❶

❷

The `join` function of `os.path` constructs a pathname out of one or more partial pathnames. In this case, it simply concatenates strings. (Note that dealing with pathnames on Windows is annoying because the backslash character must be escaped.)

❸ In this slightly less trivial case, `join` will add an extra backslash to the pathname before joining it to the filename. I was overjoyed when I discovered this, since `addSlashIfNecessary` is one of the stupid little functions I always need to write when building up my toolbox in a new language. *Do not* write this stupid little function in Python; smart people have already taken care of it for you.

❹ `expanduser` will expand a pathname that uses ~ to represent the current user's home directory. This works on any platform where users have a home directory, like Windows, UNIX, and Mac OS X; it has no effect on Mac OS.

❺ Combining these techniques, you can easily construct pathnames for directories and files under the user's home directory.

### Example 6.17. Splitting Pathnames

```
>>> os.path.split("c:\\music\\ap\\mahadeva.mp3")                    ❶
('c:\\music\\ap', 'mahadeva.mp3')
>>> (filepath, filename) = os.path.split("c:\\music\\ap\\mahadeva.mp3") ❷
>>> filepath                                                        ❸
'c:\\music\\ap'
>>> filename                                                        ❹
'mahadeva.mp3'
>>> (shortname, extension) = os.path.splitext(filename)             ❺
>>> shortname
'mahadeva'
>>> extension
'.mp3'
```

❶ The `split` function splits a full pathname and returns a tuple containing the path and filename. Remember when I said you could use multi−variable assignment to return multiple values from a function? Well, `split` is such a function.

❷ You assign the return value of the `split` function into a tuple of two variables. Each variable receives the value of the corresponding element of the returned tuple.

❸ The first variable, `filepath`, receives the value of the first element of the tuple returned from `split`, the file path.

❹ The second variable, `filename`, receives the value of the second element of the tuple returned from `split`, the filename.

❺ `os.path` also contains a function `splitext`, which splits a filename and returns a tuple containing the filename and the file extension. You use the same technique to assign each of them to separate variables.

### Example 6.18. Listing Directories

```
>>> os.listdir("c:\\music\\_singles\\")                             ❶
['a_time_long_forgotten_con.mp3', 'hellraiser.mp3',
'kairo.mp3', 'long_way_home1.mp3', 'sidewinder.mp3',
'spinning.mp3']
>>> dirname = "c:\\"
>>> os.listdir(dirname)                                             ❷
['AUTOEXEC.BAT', 'boot.ini', 'CONFIG.SYS', 'cygwin',
'docbook', 'Documents and Settings', 'Incoming', 'Inetpub', 'IO.SYS',
'MSDOS.SYS', 'Music', 'NTDETECT.COM', 'ntldr', 'pagefile.sys',
'Program Files', 'Python20', 'RECYCLER',
'System Volume Information', 'TEMP', 'WINNT']
>>> [f for f in os.listdir(dirname)
...     if os.path.isfile(os.path.join(dirname, f))]               ❸
```

```
['AUTOEXEC.BAT', 'boot.ini', 'CONFIG.SYS', 'IO.SYS', 'MSDOS.SYS',
'NTDETECT.COM', 'ntldr', 'pagefile.sys']
>>> [f for f in os.listdir(dirname)
...      if os.path.isdir(os.path.join(dirname, f))]   ❹
['cygwin', 'docbook', 'Documents and Settings', 'Incoming',
'Inetpub', 'Music', 'Program Files', 'Python20', 'RECYCLER',
'System Volume Information', 'TEMP', 'WINNT']
```

❶     The `listdir` function takes a pathname and returns a list of the contents of the directory.

❷     `listdir` returns both files and folders, with no indication of which is which.

❸     You can use list filtering and the `isfile` function of the `os.path` module to separate the files from the folders. `isfile` takes a pathname and returns 1 if the path represents a file, and 0 otherwise. Here you're using `os.path.join` to ensure a full pathname, but `isfile` also works with a partial path, relative to the current working directory. You can use `os.getcwd()` to get the current working directory.

❹     `os.path` also has a `isdir` function which returns 1 if the path represents a directory, and 0 otherwise. You can use this to get a list of the subdirectories within a directory.

**Example 6.19. Listing Directories in `fileinfo.py`**

```
def listDirectory(directory, fileExtList):
    "get list of file info objects for files of particular extensions"
    fileList = [os.path.normcase(f)
                for f in os.listdir(directory)]           ❶ ❷
    fileList = [os.path.join(directory, f)
                for f in fileList
                if os.path.splitext(f)[1] in fileExtList]  ❸ ❹ ❺
```

❶     `os.listdir(directory)` returns a list of all the files and folders in `directory`.

❷     Iterating through the list with f, you use `os.path.normcase(f)` to normalize the case according to operating system defaults. `normcase` is a useful little function that compensates for case–insensitive operating systems that think that `mahadeva.mp3` and `mahadeva.MP3` are the same file. For instance, on Windows and Mac OS, `normcase` will convert the entire filename to lowercase; on UNIX–compatible systems, it will return the filename unchanged.

❸     Iterating through the normalized list with f again, you use `os.path.splitext(f)` to split each filename into name and extension.

❹     For each file, you see if the extension is in the list of file extensions you care about (`fileExtList`, which was passed to the `listDirectory` function).

❺     For each file you care about, you use `os.path.join(directory, f)` 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` and `os.path` for file, directory, and path manipulations. These modules are wrappers for platform–specific modules, so functions like `os.path.split` work on UNIX, Windows, Mac OS, and any other platform supported by Python.

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 from working on the command line.

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

```
>>> os.listdir("c:\\music\\_singles\\")                   ❶
['a_time_long_forgotten_con.mp3', 'hellraiser.mp3',
'kairo.mp3', 'long_way_home1.mp3', 'sidewinder.mp3',
'spinning.mp3']
```

```
>>> import glob
>>> glob.glob('c:\\music\\_singles\\*.mp3')          ❷
['c:\\music\\_singles\\a_time_long_forgotten_con.mp3',
 'c:\\music\\_singles\\hellraiser.mp3',
 'c:\\music\\_singles\\kairo.mp3',
 'c:\\music\\_singles\\long_way_home1.mp3',
 'c:\\music\\_singles\\sidewinder.mp3',
 'c:\\music\\_singles\\spinning.mp3']
>>> glob.glob('c:\\music\\_singles\\s*.mp3')          ❸
['c:\\music\\_singles\\sidewinder.mp3',
 'c:\\music\\_singles\\spinning.mp3']
>>> glob.glob('c:\\music\\*\\*.mp3')                    ❹
```

❶    As you saw earlier, `os.listdir` simply takes a directory path and lists all files and directories in that directory.

❷    The `glob` 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 a directory path plus "*.mp3", which will match all `.mp3` files. Note that each element of the returned list already includes the full path of the file.

❸    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.

❹    Now consider this scenario: you have a `music` directory, with several subdirectories within it, with `.mp3` files within each subdirectory. You can get a list of all of those with a single call to `glob`, by using two wildcards at once. One wildcard is the `"*.mp3"` (to match `.mp3` files), and one wildcard is *within the directory path itself*, to match any subdirectory within `c:\music`. That's a crazy amount of

❶

❷
❸

❹
❺
❻

❶

`listDirectory` is the main attraction of this entire module. It takes a directory (like `c:\music\_singles\` in my case) and a list of interesting file extensions (like `['.mp3']`), and it returns a list of class instances that act like dictionaries that contain metadata about each interesting file in that directory. And it does it in just a few straightforward lines of code.

**❷** As you saw in the previous section, this line of code gets a list of the full pathnames of all the files in `directory` that have an interesting file extension (as specified by `fileExtList`).

**❸** Old–school Pascal programmers may be familiar with them, but most people give me a blank stare when I tell them that Python supports *nested functions* –– literally, a function within a function. The nested function `getFileInfoClass` can be called only from the function in which it is defined, `listDirectory`. As with any other function, you don't need an interface declaration or anything fancy; just define the function and code it.

**❹** Now that you've seen the `os` module, this line should make more sense. It gets the extension of the file (`os.path.splitext(filename)[1]`), forces it to uppercase (`.upper()`), slices off the dot (`[1:]`), and constructs a class name out of it with string formatting. So `c:\music\ap\mahadeva.mp3` becomes `.mp3` becomes `.MP3` becomes `MP3` becomes `MP3FileInfo`.

**❺** Having constructed the name of the handler class that would handle this file, you check to see if that handler class actually exists in this module. If it does, you return the class, otherwise you return the base class `FileInfo`. This is a very important point: *this function returns a class*. Not an instance of a class, but the class itself.

**❻** For each file in the "interesting files" list (

```
import sys
from UserDict import UserDict

def stripnulls(data):
    "strip whitespace and nulls"
    return data.replace("\00", "").strip()

class FileInfo(UserDict):
    "store file metadata"
    def __init__(self, filename=None):
        UserDict.__init__(self)
        self["name"] = filename

class MP3FileInfo(FileInfo):
    "store ID3v1.0 MP3 tags"
    tagDataMap = {"title"   : (  3,  33, stripnulls),
```

- Catching exceptions with `try...except`
- Protecting external resources with `try...finally`
- 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've used regular expressions in other languages (like Perl), the syntax will be very familiar, and you get by just reading the summary of the `re` module (http://www.python.org/doc/current/lib/module−re.html) to get an overview of the available functions and their arguments.

## 7.1. Diving In

Strings have methods for searching (`index`, `find`, and `count`), replacing (`replace`), and parsing (`split`), but they are limited to the simplest of cases. The search methods look for a single, hard−coded substring, and they are always case−sensitive. To do case−insensitive searches of a string `s`, you must call `s.lower()` or `s.upper()` and make sure your search strings are the appropriate case to match. The `replace` and `split` methods have the same limitations.

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 there's a lot to be said for fast, simple, readable code. But if you find yourself using a lot of different string functions with `if` statements to handle special cases, or if you're combining them with `split` and `join`

❶

❷

❸

❹
❺ ❻

❶

❷

`replace` method sees these two occurrences and blindly replaces both of them; meanwhile, I see my addresses getting destroyed.

❸ To solve the problem of addresses with more than one `'ROAD'` substring, you could resort to something like this: only search and replace `'ROAD'` in the last four characters of the address (`s[-4:]`), and leave the string alone (`s[:-4]`). But you can see that this is already getting unwieldy. For example, the pattern is dependent on the length of the string you're replacing (if you were replacing `'STREET'` with `'ST.'`, you would need to use `s[:-6]` and `s[-6:].replace(...)`). Would you like to come back in six months and debug this? I know I wouldn't.

❹ It's time to move up to regular expressions. In Python, all functionality related to regular expressions is contained in the `re` module.

❺ Take a look at the first parameter: `'ROAD$'`. This is a simple regular expression that matches `'ROAD'` only when it occurs at the end of a string. The `$` means "end of the string". (There is a corresponding character, the caret `^`, which means "beginning of the string".)

❻ Using the `re.sub` function, you search the string `s` for the regular expression `'ROAD$'` and replace it with `'RD.'`. This matches the ROAD at the end of the string `s`, but does *not* match the ROAD that's part of the word BROAD, because that's in the middle of `s`.

Continuing with my story of scrubbing addresses, I soon discovered that the previous example, matching `'ROAD'` at the end of the address, was not good enough, because not all addresses included a street designation at all; some just ended with the street name. Most of the time, I got away with it, but if the street name was `'BROAD'`, then the regular expression would match `'ROAD'` at the end of the string as part of the word `'BROAD'`, which is not what I wanted.

**Example 7.2. Matching Whole Words**

```
>>> s = '100 BROAD'
>>> re.sub('ROAD$', 'RD.', s)
'100 BRD.'
>>> re.sub('\\bROAD$', 'RD.', s)        ❶
'100 BROAD'
>>> re.sub(r'\bROAD$', 'RD.', s)        ❷
'100 BROAD'
>>> s = '100 BROAD ROAD APT. 3'
>>> re.sub(r'\bROAD$', 'RD.', s)        ❸
'100 BROAD ROAD APT. 3'
>>> re.sub(r'\bROAD\b', 'RD.', s)       ❹
'100 BROAD RD. APT 3'
```

❶ What I *really* wanted was to match `'ROAD'` when it was at the end of the string *and* it was its own whole word, not a part of some larger word. To express this in a regular expression, you use `\b`, which means "a word boundary must occur right here". In Python, this is complicated by the fact that the `'\'` character in a string must itself be escaped. This is sometimes referred to as the backslash plague, and it is one reason why regular expressions are easier in Perl than in Python. On the down side, Perl mixes regular expressions with other syntax, so if you have a bug, it may be hard to tell whether it's a bug in syntax or a bug in your regular expression.

❷ 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 that nothing in this string should be escaped; `'\t'` is a tab character, but `r'\t'` is really the backslash character `\` followed by the letter `t`. I recommend always using raw strings when dealing with regular expressions; otherwise, things get too confusing too quickly (and regular expressions get confusing quickly enough all by themselves).

❸ *sigh* Unfortunately, I soon found more cases that contradicted my logic. In this case, the street address contained the word `'ROAD'` as a whole word by itself, but it wasn't at the end, because the address had an apartment number after the street designation. Because `'ROAD'` isn't at the very end of

the string, it doesn't match, so the entire call to `re.sub` ends up replacing nothing at all, and you get the original string back, which is not what you want.

❹  To solve this problem, I removed the `$` character and added another `\b`. Now the regular expression reads "match `'ROAD'` when it's a whole word by itself anywhere in the string," whether at the end, the beginning, or somewhere in the middle.

# 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 ("Copyright `MCMXLVI`" instead of "Copyright `1946`"), or on the dedication walls of libraries or universities ("established `MDCCCLXXXVIII`" instead of "established `1888`"). You may also have seen them in outlines and bibliographical references. It's a system of representing numbers that really does date back to the ancient Roman empire (hence the name).

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. `I` is 1, `II` is 2, and `III` is 3. `VI` is 6 (literally, "5 and 1"), `VII` is 7, and `VIII` is 8.
- The tens characters (`I`, `X`, `C`, and `M`) can be repeated up to three times. At 4, you need to subtract from the next highest fives character. You can't represent 4 as `IIII`; instead, it is represented as `IV` ("1 less than 5"). The number 40 is written as `XL` (10 less than 50), 41 as `XLI`, 42 as `XLII`, 43 as `XLIII`, and then 44 as `XLIV` (10 less than 50, then 1 less than 5).
- Similarly, at 9, you need to subtract from the next highest tens character: 8 is `VIII`, but 9 is `IX` (1 less than 10), not `VIIII` (since the `I` character can not be repeated four times). The number 90 is `XC`, 900 is `CM`.
- The fives characters can not be repeated. The number 10 is always represented as `X`, never as `VV`. The number 100 is always `C`, never `LL`.
- 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; `CD` is a completely different number (400, 100 less than 500). `CI` is 101; `IC` is not even a valid Roman numeral (because you can't subtract 1 directly from 100; you would need to write it as `XCIX`, for 10 less than 100, then 1 less than 10).

## 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 always written highest to lowest, let's start with the highest: the thousands place. For numbers 1000 and higher, the thousands are represented by a series of `M` characters.

**Example 7.3. Checking for Thousands**

```
>>> import re
>>> pattern = '^M?M?M?$'           ❶
>>> re.search(pattern, 'M')         ❷
<SRE_Match object at 0106FB58>
```

```
>>> re.search(pattern, 'MM')        ❸
<SRE_Match object at 0106C290>
>>> re.search(pattern, 'MMM')       ❹
<SRE_Match object at 0106AA38>
>>> re.search(pattern, 'MMMM')      ❺
>>> re.search(pattern, '')          ❻
<SRE_Match object at 0106F4A8>
```

❶  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 the M characters were, which is not what you want. You want to make sure that the M characters, if they're there, are at the beginning of the string.
- M? to optionally match a single M character. Since this is repeated three times, you're matching anywhere from zero to three M characters in a row.
- $ to match what precedes only at the end of the string. When combined with the ^ character at the beginning, this means that the pattern must match the entire string, with no other characters before or after the M characters.

❷  The essence of the re module is the search function, that takes a regular expression (pattern) and a string ('M') to try to match against the regular expression. If a match is found, search returns an object which has various methods to describe the match; if no match is found, search returns None, the Python null value. All you care about at the moment is whether the pattern matches, which you can tell by just looking at the return value of search. 'M' matches this regular expression, because the first optional M matches and the second and third optional M characters are ignored.

❸  'MM' matches because the first and second optional M characters match and the third M is ignored.

❹  'MMM' matches because all three M characters match.

❺  'MMMM' does not match. All three M characters match, but then the regular expression insists on the string ending (because of the $ character), and the string doesn't end yet (because of the fourth M). So search returns None.

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 its value.

- 100 = C
- 200 = CC
- 300 = CCC
- 400 = CD
- 500 = D
- 600 = DC
- 700 = DCC
- 800 = DCCC
- 900 = CM

So there are four possible patterns:

- CM
- CD
  Zero to three C

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.

❶
❷

❸

❹

❺
❻

❶

❷

❸

❹

❺

❻

```
>>> import re
>>> pattern = '^M?M?M?$'
>>> re.search(pattern, 'M')          ❶
<_sre.SRE_Match object at 0x008EE090>
>>> pattern = '^M?M?M?$'
>>> re.search(pattern, 'MM')         ❷
<_sre.SRE_Match object at 0x008EEB48>
>>> pattern = '^M?M?M?$'
>>> re.search(pattern, 'MMM')        ❸
<_sre.SRE_Match object at 0x008EE090>
>>> re.search(pattern, 'MMMM')       ❹
>>>
```

❶  This matches the start of the string, and then the first optional M, but not the second and third M (but that's okay because they're optional), and then the end of the string.

❷  This matches the start of the string, and then the first and second optional M, but not the third M (but that's okay because it's optional), and then the end of the string.

❸  This matches the start of the string, and then all three optional M, and then the end of the string.

❹  This matches the start of the string, and then all three optional M, but then does not match the the end of the string (because there is still one unmatched M), so the pattern does not match and returns None.

**Example 7.6. The New Way: From n o m**

```
>>> pattern = '^M{0,3}$'             ❶
>>> re.search(pattern, 'M')          ❷
<_sre.SRE_Match object at 0x008EEB48>
>>> re.search(pattern, 'MM')         ❸
<_sre.SRE_Match object at 0x008EE090>
>>> re.search(pattern, 'MMM')        ❹
<_sre.SRE_Match object at 0x008EEDA8>
>>> re.search(pattern, 'MMMM')       ❺
>>>
```

❶  This pattern says: "Match the start of the string, then anywhere from zero to three M characters, then the end of the string." The 0 and 3 can be any numbers; if you want to match at least one but no more than three M characters, you could say M{1,3}.

❷  This matches the start of the string, then one M out of a possible three, then the end of the string.

❸  This matches the start of the string, then two M out of a possible three, then the end of the string.

❹  This matches the start of the string, then three M out of a possible three, then the end of the string.

❺  This matches the start of the string, then three M out of a possible three, but then *does not match* the end of the string. The regular expression allows for up to only three M characters before the end of the string, but you have four, so the pattern does not match and returns None.

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 make sure they behave the same way on all relevant inputs. You'll talk more about writing test cases later in this book.

## 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**

```
>>> pattern = '^M?M?M?M?(CM|CD|D?C?C?C?)(XC|XL|L?X?X?X?)$'
>>> re.search(pattern, 'MCMXL')        ❶
<_sre.SRE_Match object at 0x008EEB48>
>>> re.search(pattern, 'MCML')         ❷
<_sre.SRE_Match object at 0x008EEB48>
>>> re.search(pattern, 'MCMLX')        ❸
<_sre.SRE_Match object at 0x008EEB48>
>>> re.search(pattern, 'MCMLXXX')      ❹
<_sre.SRE_Match object at 0x008EEB48>
>>> re.search(pattern, 'MCMLXXXX')     ❺
>>>
```

❶   This matches the start of the string, then the first optional M, then CM, then XL, then the end of the string. Remember, the (A|B|C) syntax means "match exactly one of A, B, or C". You match XL, so you ignore the XC and L?X?X?X? choices, and then move on to the end of the string. MCML is the Roman numeral representation of 1940.

❷   This matches the start of the string, then the first optional M, then CM, then L?X?X?X?. Of the L?X?X?X?, it matches the L and skips all three optional X characters. Then you move to the end of the string. MCML is the Roman numeral representation of 1950.

❸   This matches the start of the string, then the first optional M, then CM, then the optional L and the first optional X, skips the second and third optional X, then the end of the string. MCMLX is the Roman numeral representation of 1960.

❹   This matches the start of the string, then the first optional M, then CM, then the optional L and all three optional X characters, then the end of the string. MCMLXXX is the Roman numeral representation of 1980.

❺   This matches the start of the string, then the first optional M, then CM, then the optional L and all three optional X characters, then *fails to match* the end of the string because there is still one more X unaccounted for. So the entire pattern fails to match, and returns None. MCMLXXXX is not a valid Roman numeral.

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

```
>>> pattern = '^M?M?M?M?(CM|CD|D?C?C?C?)(XC|XL|L?X?X?X?)(IX|IV|V?I?I?I?)$'
```

❶

❷

❸

❹

❶

❷

Roman numeral representation of `2666`.

❸ This matches the start of the string, then four out of four `M` characters, then `D?C{0,3}` with a `D` and three out of three `C` characters; then `L?X{0,3}` with an `L` and three out of three `X` characters; then `V?I{0,3}` with a `V` and three out of three `I` characters; then the end of the string. `MMMMDCCCLXXXVIII` is the Roman numeral representation of `3888`, and it's the longest Roman numeral you can write without extended syntax.

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 string, then zero out of four `M`, then matches `D?C{0,3}` by skipping the optional `D` and matching zero out of three `C`, then matches `L?X{0,3}` by skipping the optional `L` and matching zero out of three `X`, then matches `V?I{0,3}` by skipping the optional `V` and matching one out of three `I`. Then the end of the string. Whoa.

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 regular expressions, in the middle of a critical function of a large program. Or even imagine coming back to your own regular expressions a few months later. I've done it, and it's not a pretty sight.

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

```
<_sre.SRE_Match object at 0x008EEB48>
>>> re.search(pattern, 'M')                                    ❹
```

❶  The most important thing to remember when using verbose regular expressions is that you need to pass
    an extra argument when working with them: `re.VERBOSE` is a constant defined in the `re` module that
    signals that the pattern should be treated as a verbose regular expression. As you can see, this pattern
    has quite a bit of whitespace (all of which is ignored), and several comments (all of which are ignored).

❷

❸

❹

                                                                ❶
                                                                ❷

```
>>> phonePattern.search('800-555-1212-1234')                          ❸
>>>
```

❶ 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 {3} means "match exactly three numeric digits"; it's a variation on the {n,m} syntax you saw earlier. \d means "any numeric digit" (0 through 9). Putting it in parentheses means "match exactly three numeric digits, *and then remember them as a group that I can ask for later*". Then match a literal hyphen. Then match another group of exactly three digits. Then another literal hyphen. Then another group of exactly four digits. Then match the end of the string.

❷ To get access to the groups that the regular expression parser remembered along the way, use the groups() method on the object that the search function returns. It will return a tuple of however many groups were defined in the regular expression. In this case, you defined three groups, one with three digits, one with three digits, and one with four digits.

❸ 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 expand the regular expression.

### Example 7.11. Finding the Extension

```
>>> phonePattern = re.compile(r'^(\d{3})-(\d{3})-(\d{4})-(\d+)$')      ❶
>>> phonePattern.search('800-555-1212-1234').groups()                 ❷
('800', '555', '1212', '1234')
>>> phonePattern.search('800 555 1212 1234')                          ❸
>>>
>>> phonePattern.search('800-555-1212')                               ❹
>>>
```

❶ This regular expression is almost identical to the previous one. Just as before, you match the beginning of the string, then a remembered group of three digits, then a hyphen, then a remembered group of three digits, then a hyphen, then a remembered group of four digits. What's new is that you then match another hyphen, and a remembered group of one or more digits, then the end of the string.

❷ The groups() method now returns a tuple of four elements, since the regular expression now defines four groups to remember.

❸ Unfortunately, this regular expression is not the final answer either, because it assumes that the different parts of the phone number are separated by hyphens. What if they're separated by spaces, or commas, or dots? You need a more general solution to match several different types of separators.

❹ 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 *without* an extension. That's not what you wanted at all; if the extension is there, you want to know what it is, but if it's not there, you still want to know what the different parts of the main number are.

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

### Example 7.12. Handling Different Separators

```
>>> phonePattern = re.compile(r'^(\d{3})\D+(\d{3})\D+(\d{4})\D+(\d+)$')  ❶
>>> phonePattern.search('800 555 1212 1234').groups()                 ❷
('800', '555', '1212', '1234')
>>> phonePattern.search('800-555-1212-1234').groups()                 ❸
('800', '555', '1212', '1234')
>>> phonePattern.search('80055512121234')                             ❹
>>>
>>> phonePattern.search('800-555-1212')                               ❺
>>>
```

❶ Hang on to your hat. You're matching the beginning of the string, then a group of three digits, then \D+. What the heck is that? Well, \D matches any character *except* a numeric digit, and + means "1 or more".

❷

❸
❹

❹

❶
❷

❸

❹

❺

❶

❷

❸
❹

❺

❶
❷

❸

❹

❶ This is the same as in the previous example, except now you're matching \D*, zero or more non–numeric characters, before the first remembered group (the area code). Notice that you're not remembering these non–numeric characters (they're not in parentheses). If you find them, you'll just skip over them and then start remembering the area code whenever you get to it.

❷ 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 already handled; it's treated as a non–numeric separator and matched by the \D* after the first remembered group.)

❸ Just a sanity check to make sure you haven't broken anything that used to work. Since the leading characters are entirely optional, this matches the beginning of the string, then zero non–numeric characters, then a remembered group of three digits (800), then one non–numeric character (the hyphen), then a remembered group of three digits (555), then one non–numeric character (the hyphen), then a remembered group of four digits (1212), then zero non–numeric characters, then a remembered group of zero digits, then the end of the string.

❹ 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 there's a 1 before the area code, but you assumed that all the leading characters before the area code were non–numeric characters (\D*). Aargh.

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 be an indeterminate amount of stuff at the beginning of the string that you want to ignore. Rather than trying to match it all just so you can skip over it, let's take a different approach: don't explicitly match the beginning of the string at all. This approach is shown in the next example.


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

```
>>> phonePattern = re.compile(r'(\d{3})\D*(\d{3})\D*(\d{4})\D*(\d*)$')   ❶
>>> phonePattern.search('work 1-(800) 555.1212 #1234').groups()          ❷
('800', '555', '1212', '1234')
>>> phonePattern.search('800-555-1212')                                  ❸
('800', '555', '1212', '')
>>> phonePattern.search('80055512121234')                                ❹
('800', '555', '1212', '1234')
```

❶ Note the lack of ^ in this regular expression. You are not matching the beginning of the string anymore. There's nothing that says you need to match the entire input with your regular expression. The regular expression engine will do the hard work of figuring out where the input string starts to match, and go from there.

❷ Now you can successfully parse a phone number that includes leading characters and a leading digit, plus any number of any kind of separators around each part of the phone number.

❸ Sanity check. this still works.

❹ 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 between one and the next?

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), let's write it out as a verbose regular expression, before you forget why you made the choices you made.


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

```
>>> phonePattern = re.compile(r'''
                # don't match beginning of string, number can start anywhere
    (\d{3})     # area code is 3 digits (e.g. '800')
```

```
    \D*          # optional separator is any number of non-digits
    (\d{3})      # trunk is 3 digits (e.g. '555')
    \D*          # optional separator
    (\d{4})      # rest of number is 4 digits (e.g. '1212')
    \D*          # optional separator
    (\d*)        # extension is optional and can be any number of digits
    $            # end of string
    ''', re.VERBOSE)
>>> phonePattern.search('work 1-(800) 555.1212 #1234').groups()        ❶
('800', '555', '1212', '1234')
>>> phonePattern.search('800-555-1212')                                ❷
('800', '555', '1212', '')
```

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 parses the same inputs.

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 in Python.
- *Python Library Reference* (http://www.python.org/doc/current/lib/) summarizes the `re` module (http://www.python.org/doc/current/lib/module−re.html).

# 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 them now, believe me, you ain't seen nothing yet.

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.
- `(a|b|c)` matches either `a` or `b` or `c`.
- `(x)` in general is a *remembered group*. You can get the value of what matched by using the `groups()` method of the object returned by `re.search`.

Regular expressions are extremely powerful, but they are not the correct solution for every problem. You should learn enough about them to know when they are appropriate, when they will solve your problems, and when they will cause more problems than they solve.

> Some people, when confronted with a problem, think "I know, I'll use regular expressions."
> Now they have two problems.
> > ––Jamie Zawinski, in comp.emacs.xemacs
> > (http://groups.google.com/groups?selm=33F0C496.370D7C45%40netscape.com)

# 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|images|links] in my HTML document?" "How do I parse/translate/munge the text of my HTML document but leave the tags alone?" "How can I add/remove/quote attributes of all my HTML tags at once?" This chapter will answer all of these questions.

Here is a complete, working Python program in two parts. The first part, `BaseHTMLProcessor.py`, is a generic tool to help you process HTML files by walking through the tags and text blocks. The second part, `dialect.py`, is an example of how to use `BaseHTMLProcessor.py` to translate the text of an HTML document but leave the tags alone. Read the `doc strings` and comments to get an overview of what's going on. Most of it will seem like black magic, because it's not obvious how any of these class methods ever get called. Don't worry, all will be revealed in due time.

**Example 8.1. `BaseHTMLProcessor.py`**

If you have not already done so, you can download this and other examples (http://diveintopython.org/download/diveintopython−examples−5.4.zip) used in this book.

```
from sgmllib import SGMLParser
import htmlentitydefs

class BaseHTMLProcessor(SGMLParser):
    def reset(self):
        # extend (called by SGMLParser.__init__)
```

```
        # called for each entity reference, e.g. for "&copy;", ref will be "copy"
        # Reconstruct the original entity reference.
        self.pieces.append("&%(ref)s" % locals())
        # standard HTML entities are closed with a semicolon; other entities are not
        if htmlentitydefs.entitydefs.has_key(ref):
            self.pieces.append(";")

    def handle_data(self, text):
        # called for each block of plain text, i.e. outside of any tag and
        # not containing any character or entity references
        # Store the original text verbatim.
        self.pieces.append(text)

    def handle_comment(self, text):
        # called for each HTML comment, e.g. <!-- insert Javascript code here -->
        # Reconstruct the original comment.
        # It is especially important that the source document enclose client-side
        # code (like Javascript) within comments so it can pass through this
        # processor undisturbed; see comments in unknown_starttag for details.
        self.pieces.append("<!--%(text)s-->" % locals())

    def handle_pi(self, text):
        # called for each processing instruction, e.g. <?instruction>
        # Reconstruct original processing instruction.
        self.pieces.append("<?%(text)s>" % locals())

    def handle_decl(self, text):
        # called for the DOCTYPE, if present, e.g.
        # <!DOCTYPE html PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN"
        #     "http://www.w3.org/TR/html4/loose.dtd">
        # Reconstruct original DOCTYPE
        self.pieces.append("<!%(text)s>" % locals())

    def output(self):
        """Return processed HTML as a single string"""
        return "".join(self.pieces)
```

**Example 8.2. `dialect.py`**

```
import re
from BaseHTMLProcessor import BaseHTMLProcessor

class Dialectizer(BaseHTMLProcessor):
    subs = ()

    def reset(self):
        # extend (called from __init__ in ancestor)
        # Reset all data attributes
        self.verbatim = 0
        BaseHTMLProcessor.reset(self)

    def start_pre(self, attrs):
        # called for every <pre> tag in HTML source
        # Increment verbatim mode count, then handle tag like normal
        self.verbatim += 1
        self.unknown_starttag("pre", attrs)

    def end_pre(self):
        # called for every </pre> tag in HTML source
        # Decrement verbatim mode count
        self.unknown_endtag("pre")
```

```
            self.verbatim -= 1

    def handle_data(self, text):
        # override
        # called for every block of text in HTML source
        # If in verbatim mode, save text unaltered;
        # otherwise process the text with a series of substitutions
        self.pieces.append(self.verbatim and text or self.process(text))

    def process(self, text):
        # called from handle_data
        # Process text block by performing series of regular expression
        # substitutions (actual substitions are defined in descendant)
        for fromPattern, toPattern in self.subs:
            text = re.sub(fromPattern, toPattern, text)
        return text

class ChefDialectizer(Dialectizer):
    """convert HTML to Swedish Chef-speak

    based on the classic chef.x, copyright (c) 1992, 1993 John Hagerman
    """
    subs = ((r'a([nu])', r'u\1'),
            (r'A([nu])', r'U\1'),
            (r'a\B', r'e'),
            (r'A\B', r'E'),
            (r'en\b', r'ee'),
            (r'\Bew', r'oo'),
            (r'\Be\b', r'e-a'),
            (r'\be', r'i'),
            (r'\bE', r'I'),
            (r'\Bf', r'ff'),
            (r'\Bir', r'ur'),
            (r'(\w*?)i(\w*?)$', r'\1ee\2'),
            (r'\bow', r'oo'),
            (r'\bo', r'oo'),
            (r'\bO', r'Oo'),
            (r'the', r'zee'),
            (r'The', r'Zee'),
            (r'th\b', r't'),
            (r'\Btion', r'shun'),
            (r'\Bu', r'oo'),
            (r'\BU', r'Oo'),
            (r'v', r'f'),
            (r'V', r'F'),
            (r'w', r'w'),
            (r'W', r'W'),
            (r'([a-z])[.]', r'\1.  Bork Bork Bork!'))

class FuddDialectizer(Dialectizer):
    """convert HTML to Elmer Fudd-speak"""
    subs = ((r'[rl]', r'w'),
            (r'qu', r'qw'),
            (r'th\b', r'f'),
            (r'th', r'd'),
            (r'n[.]', r'n, uh-hah-hah-hah.'))

class OldeDialectizer(Dialectizer):
    """convert HTML to mock Middle English"""
    subs = ((r'i([bcdfghjklmnpqrstvwxyz])e\b', r'y\1'),
            (r'i([bcdfghjklmnpqrstvwxyz])e', r'y\1\1e'),
            (r'ick\b', r'yk'),
            (r'ia([bcdfghjklmnpqrstvwxyz])', r'e\1e'),
```

```python
            (r'e[ea]([bcdfghjklmnpqrstvwxyz])', r'e\1e'),
            (r'([bcdfghjklmnpqrstvwxyz])y', r'\1ee'),
            (r'([bcdfghjklmnpqrstvwxyz])er', r'\1re'),
            (r'([aeiou])re\b', r'\1r'),
            (r'ia([bcdfghjklmnpqrstvwxyz])', r'i\1e'),
            (r'tion\b', r'cioun'),
            (r'ion\b', r'ioun'),
            (r'aid', r'ayde'),
            (r'ai', r'ey'),
            (r'ay\b', r'y'),
            (r'ay', r'ey'),
            (r'ant', r'aunt'),
            (r'ea', r'ee'),
            (r'oa', r'oo'),
            (r'ue', r'e'),
            (r'oe', r'o'),
            (r'ou', r'ow'),
            (r'ow', r'ou'),
            (r'\bhe', r'hi'),
            (r've\b', r'veth'),
            (r'se\b', r'e'),
            (r"'s\b", r'es'),
            (r'ic\b', r'ick'),
            (r'ics\b', r'icc'),
            (r'ical\b', r'ick'),
            (r'tle\b', r'til'),
            (r'll\b', r'l'),
            (r'ould\b', r'olde'),
            (r'own\b', r'oune'),
            (r'un\b', r'onne'),
            (r'rry\b', r'rye'),
            (r'est\b', r'este'),
            (r'pt\b', r'pte'),
            (r'th\b', r'the'),
            (r'ch\b', r'che'),
            (r'ss\b', r'sse'),
            (r'([wybdp])\b', r'\1e'),
            (r'([rnt])\b', r'\1\1e'),
            (r'from', r'fro'),
            (r'when', r'whan'))

def translate(url, dialectName="chef"):
    """fetch URL and translate using dialect

    dialect in ("chef", "fudd", "olde")"""
    import urllib
    sock = urllib.urlopen(url)
    htmlSource = sock.read()
    sock.close()
    parserName = "%sDialectizer" % dialectName.capitalize()
    parserClass = globals()[parserName]
    parser = parserClass()
    parser.feed(htmlSource)
    parser.close()
    return parser.output()

def test(url):
```

```
        import webbrowser
        webbrowser.open_new(outfile)

if __name__ == "__main__":
    test("http://diveintopython.org/odbchelper_list.html")
```

**Example 8.3. Output of `dialect.py`**

Running this script will translate Section 3.2, Introducing Lists into mock Swedish Chef−speak (../native_data_types/chef.html) (from The Muppets), mock Elmer Fudd−speak (../native_data_types/fudd.html) (from Bugs Bunny cartoons), and mock Middle English (../native_data_types/olde.html) (loosely based on Chaucer's *The Canterbury Tales*). If you look at the HTML source of the output pages, you'll see that all the HTML tags and attributes are untouched, but the text between the tags has been "translated" into the mock language. If you look closer, you'll see that, in fact, only the titles and paragraphs were translated; the code listings and screen examples were left untouched.

```
<div class="abstract">
<p>Lists awe <span class="application">Pydon</span>'s wowkhowse datatype.
If youw onwy expewience wif wists is awways in
<span class="application">Visuaw Basic</span> ow (God fowbid) de datastowe
in <span class="application">Powewbuiwdew</span>, bwace youwsewf fow
<span class="application">Pydon</span> wists.</p>
</div>
```

## 8.2. Introducing `sgmllib.py`

HTML processing is broken into three steps: breaking down the HTML into its constituent pieces, fiddling with the pieces, and reconstructing the pieces into HTML again. The first step is done by `sgmllib.py`, a part of the standard Python library.

The key to understanding this chapter is to realize that HTML is not just text, it is structured text. The structure is derived from the more−or−less−hierarchical sequence of start tags and end tags. Usually you don't work with HTML this way; you work with it *textually* in a text editor, or *visually* in a web browser or web authoring tool. `sgmllib.py` presents HTML *structurally*.

`sgmllib.py` contains one important class: SGMLParser. SGMLParser parses HTML into useful pieces, like start tags and end tags. As soon as it succeeds in breaking down some data into a useful piece, it calls a method on itself based on what it found. In order to use the parser, you subclass the SGMLParser class and override these

*Character reference*

An escaped character referenced by its decimal or hexadecimal equivalent, like ` `. When found, `SGMLParser` calls `handle_charref` with the text of the decimal or hexadecimal character equivalent.

*Entity reference*

An HTML entity, like `&copy;`. When found, `SGMLParser` calls `handle_entityref` with the name of the HTML entity.

*Comment*

An HTML comment, enclosed in `<!-- ... -->`. When found, `SGMLParser` calls `handle_comment` with the body of the comment.

*Processing instruction*

An HTML processing instruction, enclosed in `<? ... >`. When found, `SGMLParser` calls `handle_pi` with the body of the processing instruction.

*Declaration*

An HTML declaration, such as a DOCTYPE. When found, SGMLParser vfj /F4 11 Tf.>

```
start tag: <head>
data: '\n       '
start tag: <meta http-equiv="Content-Type" content="text/html; charset=ISO-8859-1" >
data: '\n    \n       '
start tag: <title>
data: 'Dive Into Python'
end tag: </title>
data: '\n       '
start tag: <link rel="stylesheet" href="diveintopython.css" type="text/css" >
data: '\n       '

... rest of output omitted for brevity ...
```

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

- Subclass SGMLParser to create classes that extract interesting data out of HTML documents.
- Subclass SGMLParser to create BaseHTMLProcessor, which overrides all 8 handler methods and uses them to reconstruct the original HTML from the pieces.
- Subclass BaseHTMLProcessor to create Dialectizer, which adds some methods to process specific HTML tags specially, and overrides the handle_data method to provide a framework for processing the text blocks between the HTML tags.
- Subclass Dialectizer to create classes that define text processing rules used by Dialectizer.handle_data.
- Write a test suite that grabs a real web page from http://diveintopython.org/ and processes it.

Along the way, you'll also learn about locals, globals, and dictionary–based string formatting.

## 8.3. Extracting data from HTML documents

To extract data from HTML documents, subclass the SGMLParser 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, you can use file functions to read it, but the real fun begins when you get HTML from live web pages.

**Example 8.5. Introducing `urllib`**

```
>>> import urllib                                          ❶
>>> sock = urllib.urlopen("http://diveintopython.org/")    ❷
>>> htmlSource = sock.read()                               ❸
>>> sock.close()                                           ❹
>>> print htmlSource                                       ❺
<!DOCTYPE html PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN" "http://www.w3.org/TR/html4/loose.dtd":
     <meta http-equiv='Content-Type' content='text/html; charset=ISO-8859-1'>
   <title>Dive Into Python</title>
<link rel='stylesheet' href='diveintopython.css' type='text/css'>
<link rev='made' href='mailto:mark@diveintopython.org'>
<meta name='keywords' content='Python, Dive Into Python, tutorial, object-oriented, programming, docur
<meta name='description' content='a free Python tutorial for experienced programmers'>
</head>
<body bgcolor='white' text='black' link='#0000FF' vlink='#840084' alink='#0000FF'>
<table cellpadding='0' cellspacing='0' border='0' width='100%'>
<tr><td class='header' width='1%' valign='top'>diveintopython.org</td>
<td width='99%' align='right'><hr size='1' noshade></td></tr>
<tr><td class='tagline' colspan='2'>Python for experienced programmers</td></tr>
```

```
[...snip...]
```

❶  The `urllib` module is part of the standard Python library. It contains functions for getting information about

❷

❸

❹
❺

❶

❷
❸ ❹

❶

❷

❸
❹

❶
❷
❸
❹

```
        if htmlentitydefs.entitydefs.has_key(ref):
            self.pieces.append(";")

    def handle_data(self, text):              ❻
        self.pieces.append(text)

    def handle_comment(self, text):           ❼
        self.pieces.append("<!--%(text)s-->" % locals())

    def handle_pi(self, text):                ❽
        self.pieces.append("<?%(text)s" % locals())

    def handle_decl(self, text):
        self.pieces.append("<!%(text)s" % locals())
```

❶  `reset`, called by `SGMLParser.__init__`, initializes `self.pieces` as an empty list before calling the ancestor method. `self.pieces` is a data attribute which will hold the pieces of the HTML document you're constructing. Each handler method will reconstruct the HTML that `SGMLParser` parsed, and each method will append that string to `self.pieces`. Note that `self.pieces` is a list. You might be tempted to define it as a string and just keep appending each piece to it. That would work, but Python is much more efficient at dealing with lists.[2]

❷  Since `BaseHTMLProcessor` does not define any methods for specific tags (like the `start_a` method in `URLLister`), `SGMLParser` will call `unknown_starttag` for every start tag. This method takes the tag (`tag`) and the list of attribute name/value pairs (`attrs`), reconstructs the original HTML, and appends it to `self.pieces`. The string formatting here is a little strange; you'll untangle that (and also the odd–looking `locals` function) later in this chapter.

❸  Reconstructing end tags is much simpler; just take the tag name and wrap it in the `</...>` brackets.

❹  When `SGMLParser` finds a character reference, it calls `handle_charref` with the bare reference. If the HTML document contains the reference ` `, `ref` will be `160`. Reconstructing the original complete character reference just involves wrapping `ref` in `&#...;` characters.

❺  Entity references are similar to character references, but without the hash mark. Reconstructing the original entity reference requires wrapping `ref` in `&...;` characters. (Actually, as an erudite reader pointed out to me, it's slightly more complicated than this. Only certain standard HTML entites end in a semicolon; other similar–looking entities do not. Luckily for us, the set of standard HTML entities is defined in a dictionary in a Python module called `htmlentitydefs`. Hence the extra `if` statement.)

❻  Blocks of text are simply appended to `self.pieces` unaltered.

❼  HTML comments are wrapped in `<!--...-->` characters.

❽  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 do this properly (and all modern web browsers are forgiving if they don't). `BaseHTMLProcessor` is not forgiving; if script is improperly embedded, it will be parsed as if it were HTML. For instance, if the script contains less–than and equals signs, `SGMLParser` may incorrectly think that it has found tags and attributes. `SGMLParser` always converts tags and attribute names to lowercase, which may break the script, and `BaseHTMLProcessor` always encloses attribute values in double quotes (even if the original HTML document used single quotes or no quotes), which will certainly break the script. Always protect your client–side script within HTML comments.

### Example 8.9. `BaseHTMLProcessor` output

```
    def output(self):                         ❶
        """Return processed HTML as a single string"""
        return "".join(self.pieces)           ❷
```

❶

you didn't appreciate how much work Python was doing before giving you that error.

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, when you reference a variable within a nested function or `lambda` function, Python will search for that variable in the current (nested or `lambda`) function'sphat 2.v/ Tn 2le inm inle'ss. In

❶

❷

❸

❶

❷

❸

❶

variables in the local namespace.

❸ This prints `x= 1`, not `x= 2`.

❹ After being burned by `locals`, you might think that this *wouldn't* change the value of `z`, but it does. Due to internal differences in how Python is implemented (which I'd rather not go into, since I don't fully understand them myself), `globals`

❺

❶
❷
❸

❶

❷

❸

❶

❶

**Example 8.15. More dictionary–based string formatting**

```
def unknown_starttag(self, tag, attrs):
    strattrs = "".join([' %s="%s"' % (key, value) for key, value in attrs])  ❶
    self.pieces.append("<%(tag)s%(strattrs)s>" % locals())                        ❷
```

❶  When this method is called, `attrs` is a list of key/value tuples, just like the `items` of a dictionary, which means you can use multi–variable assignment to iterate through it. This should be a familiar pattern by now, but there's a lot going on here, so let's break it down:

    a. Suppose `attrs` is `[('href', 'index.html'), ('title', 'Go to home page')]`.

    b. In the first round of the list comprehension, `key` will get `'href'`, and `value` will get `'index.html'`.

    c. The string formatting `' %s="%s"' % (key, value)` will resolve to `' href="index.html"'`. This string becomes the first element of the list comprehension's return value.

    d. In the second round, `key` will get `'title'`, and `value` will get `'Go to home page'`.

    e. The string formatting will resolve to `' title="Go to home page"'`.

    f. The list comprehension returns a list of these two resolved strings, and `strattrs` will join both elements of this list together to form `' href="index.html" title="Go to home page"'`.

❷  Now, using dictionary–based string formatting, you insert the value of `tag` and `strattrs` into a string. So if `tag` is `'a'`, the final result would be `'<a href="index.html" title="Go to home page">'`, and that is what gets appended to `self.pieces`.

Using dictionary–based string formatting with `locals` is a convenient way of making complex string formatting expressions more readable, but it comes with a price. There is a slight performance hit in making the call to `locals`, since `locals` builds a copy of the local namespace.

# 8.7. Quoting attribute values,

❶

```
...          </body>
...          </html>
...          """
>>> from BaseHTMLProcessor import BaseHTMLProcessor
>>> parser = BaseHTMLProcessor()
>>> parser.feed(htmlSource)     ❷
>>> print parser.output()       ❸
<html>
<head>
<title>Test page</title>
</head>
<body>
<ul>
<li><a href="index.html">Home</a></li>
<li><a href="toc.html">Table of contents</a></li>
<li><a href="history.html">Revision history</a></li>
</body>
</html>
```

❶    Note that the attribute values of the `href` attributes in the `<a>` tags are not properly quoted. (Also note that you're using triple quotes for something other than a `doc string`. And directly in the IDE, no less. They're very useful.)

❷    Feed the parser.

❸    Using the `output` function defined in `BaseHTMLProcessor`, you get the output as a single string, complete with quoted attribute values. While this may seem anti–climactic, think about how much has actually happened here: `SGMLParser` parsed the entire HTML document, breaking it down into tags, refs, data, and so forth; `BaseHTMLProcessor` used those elements to reconstruct pieces of HTML (which are still stored in `parser.pieces`, if you want to see them); finally, you called `parser.output`, which joined all the pieces of HTML into one string.

## 8.8. Introducing `dialect.py`

`Dialectizer` is a simple (and silly) descendant of `BaseHTMLProcessor`. It runs blocks of text through a series of substitutions, but it makes sure that anything within a `<pre>...</pre>` block passes through unaltered.

To handle the `<pre>` blocks, you define two methods in `Dialectizer`: `start_pre` and `end_pre`.

**Example 8.17. Handling specific tags**

```
def start_pre(self, attrs):                  ❶
    self.verbatim += 1                       ❷
    self.unknown_starttag("pre", attrs)      ❸

def end_pre(self):                           ❹
    self.unknown_endtag("pre")               ❺
    self.verbatim -= 1                       ❻
```

❶    `start_pre` is called every time `SGMLParser` finds a `<pre>` tag in the HTML source. (In a minute, you'll see exactly how this happens.) The method takes a single parameter, `attrs`, which contains the attributes of the tag (if any). `attrs` is a list of key/value tuples, just like `unknown_starttag` takes.

❷    In the `reset` method, you initialize a data attribute that serves as a counter for `<pre>` tags. Every time you hit a `<pre>` tag, you increment the counter; every time you hit a `</pre>` tag, you'll decrement the counter. (You could just use this as a flag and set it to `1` and reset it to `0`, but it's just as easy to do it this way, and this handles the odd (but possible) case of nested `<pre>` tags.) In a minute, you'll see how this counter is put to good use.

❸

That's it, that's the only special processing you do for `<pre>` tags. Now you pass the list of attributes along to `unknown_starttag` so it can do the default processing.

❹ `end_pre` is called every time `SGMLParser` finds a `</pre>` tag. Since end tags can not contain attributes, the method takes no parameters.

❺ First, you want to do the default processing, just like any other end tag.

❻ Second, you decrement your counter to signal that this `<pre>` block has been closed.

At this point, it's worth digging a little further into `SGMLParser`. I've claimed repeatedly (and you've taken it on faith so far) that `SGMLParser` looks for and calls specific methods for each tag, if they exist. For instance, you just saw the definition of `start_pre` and `end_pre` to handle `<pre>` and `</pre>`. But how does this happen? Well, it's not magic, it's just good Python coding.

**Example 8.18. `SGMLParser`**

```
def finish_starttag(self, tag, attrs):                  ❶
    try:
        method = getattr(self, 'start_' + tag)          ❷
    except AttributeError:                              ❸
        try:
            method = getattr(self, 'do_' + tag)         ❹
        except AttributeError:
            self.unknown_starttag(tag, attrs)           ❺
            return -1
        else:
            self.handle_starttag(tag, method, attrs)    ❻
            return 0
    else:
        self.stack.append(tag)
        self.handle_starttag(tag, method, attrs)
        return 1                                        ❼

def handle_starttag(self, tag, method, attrs):
    method(attrs)                                       ❽
```

❶ At this point, `SGMLParser` 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 handler method for this tag, or whether you should fall back on the default method (`unknown_starttag`).

❷ The "magic" of `SGMLParser` is nothing more than your old friend, `getattr`. What you may not have realized before is that `getattr` will find methods defined in descendants of an object as well as the object itself. Here the object is `self`, the current instance. So if `tag` is `'pre'`, this call to `getattr` will look for a `start_pre` method on the current instance, which is an instance of the `Dialectizer` class.

❸ `getattr` 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 to `getattr` inside a `try...except` block and explicitly caught the `AttributeError`.

❹ Since you didn't find a `start_xxx` method, you'll also look for a `do_xxx` method before giving up. This alternate naming scheme is generally used for standalone tags, like `<br>`, which have no corresponding end tag. But you can use either naming scheme; as you can see, `SGMLParser` tries both for every tag. (You shouldn't define both a `start_xxx` and `do_xxx` handler method for the same tag, though; only the `start_xxx` method will get called.)

❺ Another `AttributeError`, which means that the call to `getattr` failed with `do_xxx`. Since you found neither a `start_xxx` nor a `do_xxx` method for this tag, you catch the

exception and fall back on the default method, `unknown_starttag`.

**❻** Remember, `try...except` blocks can have an `else` clause, which is called if no exception is raised during the `try...except` block. Logically, that means that you *did* find a `do_xxx` method for this tag, so you're going to call it.

**❼** 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 the `self.stack.append(tag)` either; `SGMLParser` keeps track internally of whether your start tags are balanced by appropriate end tags, but it doesn't do anything with this information either. In theory, you could use this module to validate that your tags were fully balanced, but it's probably not worth it, and it's beyond the scope of this chapter. You have better things to worry about right now.

**❽** `start_xxx` and `do_xxx` methods are not called directly; the tag, method, and attributes are passed to this function, `handle_starttag`

**❶**
**❷**

**❶**
**❷**

**❶**
**❷**

dialect, you would simply add an appropriately–named file in the plug–ins directory (like `foodialect.py` which contains the `FooDialectizer` class). Calling the `translate` function with the dialect name `'foo'` would find the module `foodialect.py`, import the class `FooDialectizer`, and away you go.

**Example 8.22. The `translate` function, part 3**

```
parser.feed(htmlSource)     ❶
parser.close()              ❷
return parser.output()      ❸
```

❶    After all that imagining, this is going to seem pretty boring, but the `feed` function is what does the entire transformation. You had the entire HTML source in a single string, so you only had to call `feed` once. However, you can call `feed` as often as you want, and the parser will just keep parsing. So if you were worried about memory usage (or you knew you were going to be dealing with very large HTML pages), you could set this up in a loop, where you read a few bytes of HTML and fed it to the parser. The result would be the same.

❷    Because `feed` 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 you did). Otherwise you may find that your output is missing the last few bytes.

❸    Remember, `output` is the function you defined on `BaseHTMLProcessor` 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/dialect/). Unfortunately, source code does not appear to be available.

# 8.10. Summary

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

- 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 `Dialectizer`

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

- Using `locals()` and `globals()` to access namespaces
- Formatting strings using dictionary–based substitutions

---

[1] The technical term for a parser like `SGMLParser` is a *consumer*: it consumes HTML and breaks it down. Presumably, the name `feed` was chosen to fit into the whole "consumer" motif. Personally, it makes me think of an exhibit in the zoo where there's just a dark cage with no trees or plants or evidence of life of any kind, but if you stand perfectly still and look really closely you can make out two beady eyes staring back at you from the far left corner, but you convince yourself that that's just your mind playing tricks on you, and the only way you can tell that the whole thing isn't just an empty cage is a small innocuous sign on the railing that reads, "Do not feed the parser." But maybe

that's just me. In any event, it's an interesting mental image.

[2] The 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 the element and updates the index. Since strings can not be changed after they are created, code like `s = s + newpiece` will create an entirely new string out of the concatenation of the original and the new piece, then throw away the original string. This involves a lot of expensive memory management, and the amount of effort involved increases as the string gets longer, so doing `s = s + newpiece` in a loop is deadly. In technical terms, appending `n` items to a list is `O(n)`, while appending `n` items to a string is `O(n`$^2$`)`.

[3] I don't get out much.

[4] All 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 better or worse than Perl?" (answer: "Perl is worse than Python because people wanted it worse." –Larry Wall, 10/14/1998) But questions about HTML processing pop up in one form or another about once a month, and among those questions, this is a popular one.

```
import getopt

_debug = 0

class NoSourceError(Exception): pass

class KantGenerator:
    """generates mock philosophy based on a context-free grammar"""

    def __init__(self, grammar, source=None):
        self.loadGrammar(grammar)
        self.loadSource(source and source or self.getDefaultSource())
        self.refresh()

    def _load(self, source):
        """load source):
```

```
def refresh(self):
    """reset output buffer, re-parse entire source file, and return output

    Since parsing involves a good deal of randomness, this is an
    easy way to get new output without having to reload a grammar file
    each time.
    """
    self.reset()
    self.parse(self.source)
    return self.output()

def output(self):
    """output generated text"""
    return "".join(self.pieces)

def randomChildElement(self, node):
    """choose a random child element of a node

    This is a utility method used by do_xref and do_choice.
    """
    choices = [e for e in node.childNodes
                 if e.nodeType == e.ELEMENT_NODE]
    chosen = random.choice(choices)
    if _debug:
        sys.stderr.write('%s available choices: %s\n' % \
            (len(choices), [e.toxml() for e in choices]))
        sys.stderr.write('Chosen: %s\n' % chosen.toxml())
    return chosen

def parse(self, node):
    """parse a single XML node

    A parsed XML document (from minidom.parse) is a tree of nodes
    of various types.  Each node is represented by an instance of the
    corresponding Python class (Element for a tag, Text for
    text data, Document for the top-level document).  The following
    statement constructs the name of a class method based on the type
    of node we're parsing ("parse_Element" for an Element node,
    "parse_Text" for a Text node, etc.) and then calls the method.
    """
    parseMethod = getattr(self, "parse_%s" % node.__class__.__name__)
    parseMethod(node)

def parse_Document(self, node):
    """parse the document node

    The document node by itself isn't interesting (to us), but
    its only child, node.documentElement, is: it's the root node
    of the grammar.
    """
    self.parse(node.documentElement)

def parse_Text(self, node):
    """parse a text node

    The text of a text node is usually added to the output buffer
    verbatim.  The one exception is that <p class='sentence'> sets
    a flag to capitalize the first letter of the next word.  If
    that flag is set, we capitalize the text and reset the flag.
    """
    text = node.data
    if self.capitalizeNextWord:
        self.pieces.append(text[0].upper())
```

```
            self.pieces.append(text[1:])
            self.capitalizeNextWord = 0
        else:
            self.pieces.append(text)

    def parse_Element(self, node):
        """parse an element

        An XML element corresponds to an actual tag in the source:
        <xref id='...'>, <p chance='...'>, <choice>, etc.
        Each element type is handled in its own method.  Like we did in
        parse(), we construct a method name based on the name of the
        element ("do_xref" for an <xref> tag, etc.) and
        call the method.
        """
        handlerMethod = getattr(self, "do_%s" % node.tagName)
        handlerMethod(node)

    def parse_Comment(self, node):
        """parse a comment

        The grammar can contain XML comments, but we ignore them
        """
        pass

    def do_xref(self, node):
        """handle <xref id='...'> tag

        An <xref id='...'> tag is a cross-reference to a <ref id='...'>
        tag.  <xref id='sentence'/> evaluates to a randomly chosen child of
        <ref id='sentence'>.
        """
        id = node.attributes["id"].value
        self.parse(self.randomChildElement(self.refs[id]))

    def do_p(self, node):
        """handle <p> tag

        The <p> tag is the core of the grammar.  It can contain almost
        anything: freeform text, <choice> tags, <xref> tags, even other
        <p> tags.  If a "class='sentence'" attribute is found, a flag
        is set and the next word will be capitalized.  If a "chance='X'"
        attribute is found, there is an X% chance that the tag will be
        evaluated (and therefore a (100-X)% chance that it will be
        completely ignored)
        """
        keys = node.attributes.keys()
        if "class" in keys:
            if node.attributes["class"].value == "sentence":
                self.capitalizeNextWord = 1
        if "chance" in keys:
            chance = int(node.attributes["chance"].value)
            doit = (chance > random.randrange(100))
        else:
            doit = 1
        if doit:
            for child in node.childNodes: self.parse(child)

    def do_choice(self, node):
        """handle <choice> tag

        A <choice> tag contains one or more <p> tags.  One <p> tag
        is chosen at random and evaluated; the rest are ignored.
```

```
            """
            self.parse(self.randomChildElement(node))

def usage():
    print __doc__

def main(argv):
    grammar = "kant.xml"
    try:
        opts, args = getopt.getopt(argv, "hg:d", ["help", "grammar="])
    except getopt.GetoptError:
        usage()
        sys.exit(2)
    for opt, arg in opts:
        if opt in ("-h", "--help"):
            usage()
            sys.exit()
        elif opt == '-d':
            global _debug
            _debug = 1
        elif opt in ("-g", "--grammar"):
            grammar = arg

    source = "".join(args)

    k = KantGenerator(grammar, source)
    print k.output()

if __name__ == "__main__":
    main(sys.argv[1:])
```

### Example 9.2. `toolbox.py`

```
"""Miscellaneous utility functions"""

def openAnything(source):
    """URI, filename, or string --> stream

    This function lets you define parsers that take any input source
    (URL, pathname to local or network file, or actual data as a string)
    and deal with it in a uniform manner.  Returned object is guaranteed
    to have all the basic stdio read methods (read, readline, readlines).
    Just .close() the object when you're done with it.

    Examples:
    >>> from xml.dom import minidom
    >>> sock = openAnything("http://localhost/kant.xml")
    >>> doc = minidom.parse(sock)
    >>> sock.close()
    >>> sock = openAnything("c:\\inetpub\\wwwroot\\kant.xml")
    >>> doc = minidom.parse(sock)
    >>> sock.close()
    >>> sock = openAnything("<ref id='conjunction'><text>and</text><text>or</text></ref>")
    >>> doc = minidom.parse(sock)
    >>> sock.close()
    """
    if hasattr(source, "read"):
        return source

    if source == '-':
        import sys
```

```
        return sys.stdin

    # try to open with urllib (if source is http, ftp, or file URL)
    import urllib
    try:
        return urllib.urlopen(source)
    except (IOError, OSError):
        pass

    # try to open with native open function (if source is pathname)
    try:
        return open(source)
    except (IOError, OSError):
        pass

    # treat source as string
    import StringIO
    return StringIO.StringIO(str(source))
```

Run the program `kgp.py` by itself, and it will parse the default XML–based grammar, in `kant.xml`, and print several paragraphs worth of philosophy in the style of Immanuel Kant.**kgp.py**

This is, of course, complete gibberish. Well, not complete gibberish. It is syntactically and grammatically correct (although very verbose –– Kant wasn't what you would call a get–to–the–point kind of guy). Some of it may actually be true (or at least the sort of thing that Kant would have agreed with), some of it is blatantly false, and most of it is simply incoherent. But all of it is in the style of Immanuel Kant.

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

The interesting thing about this program is that there is nothing Kant–specific about it. All the content in the previous example was derived from the grammar file, `kant.xml`. If you tell the program to use a different grammar file (which you can specify on the command line), the output will be completely different.

**Example 9.4. Simpler output from `kgp.py`**

```
[you@localhost kgp]$ python kgp.py -g binary.xml
00101001
[you@localhost kgp]$ python kgp.py -g binary.xml
10110100
```

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 defines the structure of the output, and the `kgp.py` program reads through the grammar and makes random decisions about which words to plug in where.

# 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 to take a short detour to talk about packages.

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

```
>>> from xml.dom import minidom ❶
>>> xmldoc = minidom.parse('~/diveintopython/common/py/kgp/binary.xml')
```

❶     This is a syntax you haven't seen before. It looks almost like the `from module import` you know and love, but the "`.`" gives it away as something above and beyond a simple import. In fact, `xml` is what is

❶

have been unwieldy (as of this writing, the XML package has over 3000 lines of code) and difficult to manage (separate source files mean multiple people can work on different areas simultaneously).

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), invest some time designing a good package architecture. It's one of the many things Python is good at, so take advantage of it.

# 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)**

```
>>> from xml.dom import minidom                                         ❶
>>> xmldoc = minidom.parse('~/diveintopython/common/py/kgp/binary.xml') ❷
>>> xmldoc                                                              ❸
<xml.dom.minidom.Document instance at 010BE87C>
>>> print xmldoc.toxml()                                                ❹
<?xml version="1.0" ?>
<grammar>
<ref id="bit">
  <p>0</p>
  <p>1</p>
</ref>
<ref id="byte">
  <p><xref id="bit"/><xref id="bit"/><xref id="bit"/><xref id="bit"/>\
<xref id="bit"/><xref id="bit"/><xref id="bit"/><xref id="bit"/></p>
</ref>
</grammar>
```

❶    As you saw in the previous section, this imports the `minidom` module from the `xml.dom` package.

❷    Here is the one line of code that does all the work: `minidom.parse` takes one argument and returns a parsed representation of the XML document. The argument can be many things; in this case, it's simply a filename of an XML document on my local disk. (To follow along, you'll need to change the path to point to your downloaded examples directory.) But you can also pass a file object, or even a file–like object. You'll take advantage of this flexibility later in this chapter.

❸    The object returned from `minidom.parse` is a `Document` object, a descendant of the `Node` class. This `Document` object is the root level of a complex tree–like structure of interlocking Python objects that completely represent the XML document you passed to `minidom.parse`.

❹    `toxml` is a method of the `Node` class (and is therefore available on the `Document` object you got from `minidom.parse`). `toxml` prints out the XML that this `Node` represents. For the `Document` node, this prints out the entire XML document.

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

**Example 9.9. Getting child nodes**

```
>>> xmldoc.childNodes        ❶
[<DOM Element: grammar at 17538908>]
>>> xmldoc.childNodes[0]     ❷
<DOM Element: grammar at 17538908>
>>> xmldoc.firstChild        ❸
<DOM Element: grammar at 17538908>
```

**❶** Every `Node` has a `childNodes` attribute, which is a list of the `Node` objects. A `Document` always has only one child node, the root element of the XML document (in this case, the `grammar` element).

**❷** 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; this is just a regular Python list of regular Python objects.

**❸** Since getting the first child node of a node is a useful and common activity, the `Node` class has a `firstChild` attribute, which is synonymous with `childNodes[0]`. (There is also a `lastChild` attribute, which is synonymous with `childNodes[-1]`.)

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

```
>>> grammarNode = xmldoc.firstChild
>>> print grammarNode.toxml()    ❶
<grammar>
<ref id="bit">
  <p>0</p>
  <p>1</p>
</ref>
<ref id="byte">
  <p><xref id="bit"/><xref id="bit"/><xref id="bit"/><xref id="bit"/>\
</a
```

**❶**

**❶**

**❷**

**❸**

**❹**

**❺**

**❶**

**❷**

❸ The second child is an `Element` object representing the first `ref` element.

❹ The fourth child is an `Element` object representing the second `ref` element.

❺ The last child is a `Text` object representing the carriage return after the `'</ref>'` end tag and before the `'</grammar>'` end tag.

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

```
>>> grammarNode
<DOM Element: grammar at 19167148>
>>> refNode = grammarNode.childNodes[1]  ❶
>>> refNode
<DOM Element: ref at 17987740>
>>> refNode.childNodes                    ❷
[<DOM Text node "\n">, <DOM Text node "  ">, <DOM Element: p at 19315844>, \
<DOM Text node "\n">, <DOM Text node "  ">, \
<DOM Element: p at 19462036>, <DOM Text node "\n">]
>>> pNode = refNode.childNodes[2]
>>> pNode
<DOM Element: p at 19315844>
>>> print pNode.toxml(]h36>,1.063 Td(>>>❸Tj (print pNode.toxml(]h36>te "  "u>0</pDOM Text node "\n">
>>> 00 cmta"  ">, <DOMThe fC 2]
                                          ❹

                                          ❺
```

❶

❷

❸
❹

❺

multilingual documents, with characters from multiple languages in the same document. (They typically used escape codes to switch modes; poof, you're in Russian koi8–r mode, so character 241 means this; poof, now you're in Mac

❶

❷

❶

❷

❶
❷

❸

of the characters in the original unicode string.

❺    Printing the `koi8-r`−encoded string will probably show gibberish on your screen, because your Python IDE is interpreting those characters as `iso-8859-1`, not `koi8-r`. But at least they do print. (And, if you look carefully, it's the same gibberish that you saw when you opened the original XML document in a non−unicode−aware text editor. Python converted it from `koi8-r` into unicode when it parsed the XML document, and you've just converted it back.)

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 XML documents are all 7−bit ASCII (like the examples in this chapter), you will literally never think about unicode. Python will convert the ASCII data in the XML documents into unicode while parsing, and auto−coerce it back to ASCII whenever necessary, and you'll never even notice. But if you need to deal with that in other languages, Python is ready.

### 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/standard/principles.html).
- Unicode Tutorial (http://www.reportlab.com/i18n/python_unicode_tutorial.html) has some more examples of how to use Python's unicode functions, including how to force Python to coerce unicode into ASCII even when it doesn't really want to.
- PEP 263 (http://www.python.org/peps/pep−0263.html) goes into more detail about how and when to define a character encoding in your `.py` files.

## 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 your XML document, there is a shortcut you can use to find it quickly: `getElementsByTagName`.

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

**Example 9.20. `binary.xml`**

```
<?xml version="1.0"?>
<!DOCTYPE grammar PUBLIC "-//diveintopython.org//DTD Kant Generator Pro v1.0//EN" "kgp.dtd">
<grammar>
<ref id="bit">
  <p>0</p>
  <p>1</p>
</ref>
<ref id="byte">
  <p><xref id="bit"/><xref id="bit"/><xref id="bit"/><xref id="bit"/>\
<xref id="bit"/><xref id="bit"/><xref id="bit"/><xref id="bit"/></p>
</ref>
</grammar>
```

It has two `ref`s, `'bit'` and `'byte'`. A `bit` is either a `'0'` or `'1'`, and a `byte` is 8 bits.

**Example 9.21. Introducing `getElementsByTagName`**

```
>>> from xml.dom import minidom
>>> xmldoc = minidom.parse('binary.xml')
>>> reflist = xmldoc.getElementsByTagName('ref')  ❶
```

```
>>> reflist
[<DOM Element: ref at 136138108>, <DOM Element: ref at 136144292>]
>>> print reflist[0].toxml()
<ref id="bit">
  <p>0</p>
  <p>1</p>
</ref>
>>> print reflist[1].toxml()
<ref id="byte">
  <p><xref id="bit"/><xref id="bit"/><xref id="bit"/><xref id="bit"/>\
<xref id="bit"/><xref id="bit"/><xref id="bit"/><xref id="bit"/></p>
</ref>
```

❶    `getElementsByTagName` 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 that name. In this case, you find two `ref` elements.

**Example 9.22. Every element is searchable**

```
>>> firstref = reflist[0]                              ❶
>>> print firstref.toxml()
<ref id="bit">
  <p>0</p>
  <p>1</p>
</ref>
>>> plist = firstref.getElementsByTagName("p")   ❷
>>> plist
[<DOM Element: p at 136140116>, <DOM Element: p at 136142172>]
>>> print plist[0].toxml()                             ❸
<p>0</p>
>>> print plist[1].toxml()
<p>1</p>
```

❶    Continuing from the previous example, the first object in your `reflist` is the `'bit'` ref element.

❷    You can use the same `getElementsByTagName` method on this `Element` to find all the `<p>` elements within the `'bit'` ref element.

❸    Just as before, the `getElementsByTagName` 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**

```
>>> plist = xmldoc.getElementsByTagName("p")   ❶
>>> plist
[<DOM Element: p at 136140116>, <DOM Element: p at 136142172>, <DOM Element: p at 136146124>]
>>> plist[0].toxml()                                   ❷
'<p>0</p>'
>>> plist[1].toxml()
'<p>1</p>'
>>> plist[2].toxml()                                   ❸
'<p><xref id="bit"/><xref id="bit"/><xref id="bit"/><xref id="bit"/>\
<xref id="bit"/><xref id="bit"/><xref id="bit"/><xref id="bit"/></p>'
```

❶    Note carefully the difference between this and the previous example. Previously, you were searching for `p` elements within `firstref`, but here you are searching for `p` elements within `xmldoc`, the root–level object that represents the entire XML document. This *does* find the `p` elements nested within the `ref` elements within the root `grammar` element.

❷    The first two `p` elements are within the first `ref` (the `'bit'` ref).

❸  The last `p` element is the one within the second `ref` (the `'byte'` ref).

## 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 `binary.xml` grammar 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 also have attributes. When you parse an XML document, you get a bunch of Python objects that represent all the pieces of the XML document, and some of these Python objects represent attributes of the XML elements. But the (Python) objects that represent the (XML) attributes also have (Python) attributes, which are used to access various parts of the (XML) attribute that the object represents. I told you it was confusing. I am open to suggestions on how to distinguish these more clearly.

**Example 9.24. Accessing element attributes**

```
>>> xmldoc = minidom.parse('binary.xml')
>>> reflist = xmldoc.getElementsByTagName('ref')
>>> bitref = reflist[0]
>>> print bitref.toxml()
<ref id="bit">
  <p>0</p>
  <p>1</p>
</ref>
>>> bitref.attributes                                    ❶
<xml.dom.minidom.NamedNodeMap instance at 0x81e0c9c>
>>> bitref.attributes.keysrg 0 ❷❸063 Td(>>> )Tj (bitref.attPu'id'y) attribute0 0.00 rg 0 -11.063 Td
                                    ❹
                                    ❺
```

❶

❷

❸
❹

❺

```
>>> a
<xml.dom.minidom.Attr instance at 0x81d5044>
>>> a.name   ❶
u'id'
>>> a.value  ❷
u'bit'
```

❶    The `Attr` object completely represents a single XML attribute of a single XML element. The
name of the attribute (the same name as you used to find this object in the

❷

# 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. But they don't. Instead, they take a *file−like object*.

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 `size` parameter, it reads everything there is to read from the input source and returns all the data as a single string. When called with a `size` parameter, it reads that much from the input source and returns that much data; when called again, it picks up where it left off and returns the next chunk of data.

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: a file on disk, a web page, even a hard−coded string. As long as you pass a file−like object to the function, and the function simply calls the object's `read` method, the function can handle any kind of input source without specific code to handle each kind.

In case you were wondering how this relates to XML processing, `minidom.parse` is one such function which can take a file−like object.

**Example 10.1. Parsing XML from a file**

```
>>> from xml.dom import minidom
>>> fsock = open('binary.xml')          ❶
>>> xmldoc = minidom.parse(fsock)        ❷
>>> fsock.close()                        ❸
>>> print xmldoc.toxml()                 ❹
<?xml version="1.0" ?>
<grammar>
<ref id="bit">
  <p>0</p>
  <p>1</p>
</ref>
<ref id="byte">
  <p><xref id="bit"/><xref id="bit"/><xref id="bit"/><xref id="bit"/>\
<xref id="bit"/><xref id="bit"/><xref id="bit"/><xref id="bit"/></p>
</ref>
</grammar>
```

❶     First, you open the file on disk. This gives you a file object.

❷     You pass the file object to `minidom.parse`, which calls the Tj 0.0Quch function which can toc umet xrom the iile−on disk,

❸

        ,witls ot ldoj 0. −13.2 Td(this ror rou )Tj 0 −26.276 Td(YCllsng Xhe fj 0.0Quc )Tj .poxml(),wan himply j 0. −13.2 Td(th

❹     tgong Xho e aarseng Xnlonallfile− you oan hass ahe file ame,0l.f j 0.0Quc

**Example 10.2. Parsing XML from a URL**

```
>>> import urllib
>>> usock = urllib.urlopen('http://slashdot.org/slashdot.rdf')   ❶
>>> xmldoc = minidom.parse(usock)                                ❷
>>> usock.close()                                                ❸
>>> print xmldoc.toxml()                                         ❹
<?xml version="1.0" ?>
<rdf:RDF xmlns="http://my.netscape.com/rdf/simple/0.9/"
 xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#">

<channel>
<title>Slashdot</title>
<link>http://slashdot.org/</link>
<description>News for nerds, stuff that matters</description>
</channel>

<image>
<title>Slashdot</title>
<url>http://images.slashdot.org/topics/topicslashdot.gif</url>
<link>http://slashdot.org/</link>
</image>

<item>
<title>To HDTV or Not to HDTV?</title>
<link>http://slashdot.org/article.pl?sid=01/12/28/0421241</link>
</item>

[...snip...]
```

❶ As you saw in a previous chapter, `urlopen` takes a web page URL and returns a file–like object. Most importantly, this object has a `read` method which returns the HTML source of the web page.

❷ Now you pass the file–like object to `minidom.parse`, which obediently calls the `read` method of the object and parses the XML data that the `read` method returns. The fact that this XML data is now coming straight from a web page is completely irrelevant. `minidom.parse` doesn't know about web pages, and it doesn't care about web pages; it just knows about file–like objects.

❸ As soon as you're done with it, be sure to close the file–like object that `urlopen` gives you.

❹ 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/), a technical news and gossip site.

**Example 10.3. Parsing XML from a string (the easy but inflexible way)**

```
>>> contents = "<grammar><ref id='bit'><p>0</p><p>1</p></ref></grammar>"
>>> xmldoc = minidom.parseString(contents)   ❶
>>> print xmldoc.toxml()
<?xml version="1.0" ?>
<grammar><ref id="bit"><p>0</p><p>1</p></ref></grammar>
```

❶ `minidom` has a method, `parseString`, which takes an entire XML document as a string and parses it. You can use this instead of `minidom.parse` if you know you already have your entire XML document in a string.

OK, so you can use the `minidom.parse` function for parsing both local files and remote URLs, but for parsing strings, you use... a different function. That means that if you want to be able to take input from a file, a URL, or a string, you'll need special logic to check whether it's a string, and call the `parseString` function instead. How unsatisfying.

If there were a way to turn a string into a file–like object, then you could simply pass this object to `minidom.parse`. And in fact, there is a module specifically designed for doing just that: `StringIO`.

Dive Into Python

### Example 10.4. Introducing `StringIO`

```
>>> contents = "<grammar><ref id='bit'><p>0</p><p>1</p></ref></grammar>"
>>> import StringIO
>>> ssock = StringIO.StringIO(contents)      ❶
>>> ssock.read()                             ❷
"<grammar><ref id='bit'><p>0</p><p>1</p></ref></grammar>"
>>> ssock.read()                             ❸
''
>>> ssock.seek(0)                            ❹
>>> ssock.read(15)                           ❺
'<grammar><ref i'
>>> ssock.read(15)
"d='bit'><p>0</p"
>>> ssock.read()
'><p>1</p></ref></grammar>'
>>> ssock.close()                            ❻
```

❶    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 when creating an instance.

❷    Now you have a file–like object, and you can do all sorts of file–like things with it. Like `readr`ample 10.4. Intr

❸

❹

❺
❻

❶

❶

❶

```
        return urllib.urlopen(source)        ❷
    except (IOError, OSError):
        pass

    # try to open with native open function (if source is pathname)
    try:
        return open(source)                   ❸
    except (IOError, OSError):
        pass

    # treat source as string
                                              ❹
```

❶

❷

❸

❹

IDE, `stdout` and `stderr` default to your "Interactive Window".)

## Example 10.8. Introducing `stdout` and `stderr`

```
>>> for i in range(3):
...     print 'Dive in'          ❶
Dive in
Dive in
Dive in
>>> import sys
>>> for i in range(3):
...     sys.stdout.write('Dive in') ❷
Dive inDive inDive in
>>> for i in range(3):
...     sys.stderr.write('Dive in') ❸
Dive inDive inDive in
```

❶     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 something a set number of times.

❷     `stdout` 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 a carriage return to the end of the string you're printing, and calls `sys.stdout.write`.

❸     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 from the command line). Like `stdout`, `stderr` does not add carriage returns for you; if you want them, add them yourself.

`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 no `read` method, only `write`. Still, they are file–like objects, and you can assign any other file– or file–like object to them to redirect their output.

## Example 10.9. Redirecting output

```
[you@localhost kgp]$ python stdout.py
Dive in
[you@localhost kgp]$ cat out.log
This message will be logged instead of displayed
```

(On Windows, you can use `type` instead of `cat` 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–5.4.zip) used in this book.

```
#stdout.py
import sys

print 'Dive in'                                          ❶
saveout = sys.stdout                                     ❷
fsock = open('out.log', 'w')                             ❸
sys.stdout = fsock                                       ❹
print 'This message will be logged instead of displayed' ❺
sys.stdout = saveout                                     ❻
fsock.close()                                            ❼
```

❶     This will print to the IDE "Interactive Window" (or the terminal, if running the script from the command line).

❷     Always save `stdout` before redirecting it, so you can set it back to normal later.

❸     Open a file for writing. If the file doesn't exist, it will be created. If the file does exist, it will be overwritten.

❹     Redirect all further output to the new file you just opened.

❺     This will be "printed" to the log file only; it will not be visible in the IDE window or on the screen.

❻     Set `stdout` back to the way it was before you mucked with it.

❼     Close the log file.

Redirecting `stderr` works exactly the same way, using `sys.stderr` instead of `sys.stdout`.

**Example 10.10. Redirecting error information**

```
[you@localhost kgp]$ python stderr.py
[you@localhost kgp]$ cat error.log
Traceback (most recent line last):
  File "stderr.py", line 5, in ?
    raise Exception, 'this error will be logged'
Exception: this error will be logged
```

If you have not already done so, you can download this and other examples
(http://diveintopython.org/download/diveintopython–examples–5.4.zip) used in this book.

```
#stderr.py
import sys

fsock = open('error.log', 'w')                    ❶
sys.stderr = fsock                                ❷
raise Exception, 'this error will be logged'  ❸ ❹
```

❶     Open the log file where you want to store debugging information.

❷     Redirect standard error by assigning the file object of the newly–opened log file to `stderr`.

❸     Raise an exception. Note from the screen output that this does *not* print anything on screen. All the normal traceback information has been written to `error.log`.

❹     Also note that you're not explicitly closing your log file, nor are you setting `stderr` back to its original value. This is fine, since once the program crashes (because of the exception), Python will clean up and close the file for us, and it doesn't make any difference that `stderr` is never restored, since, as I mentioned, the program crashes and Python ends. Restoring the original is more important for `stdout`, if you expect to go do other stuff within the same script afterwards.

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 of redirecting it outright.

**Example 10.11. Printing to `stderr`**

```
>>> print 'entering function'
entering function
>>> import sys
>>> print >> sys.stderr, 'entering function'  ❶
entering function
```

❶     This shorthand syntax of the `print` statement can be used to write to any open file, or file–like object. In this case, you can redirect a single `print` statement to `stderr` without affecting subsequent `print` statements.

Standard input, on the other hand, is a read–only file object, and it represents the data flowing into the program from some previous program. This will likely not make much sense to classic Mac OS users, or even Windows users unless you were ever fluent on the MS–DOS command line. The way it works is that you can construct a chain of commands in a single line, so that one program's output becomes the input for the next program in the chain. The first program simply outputs to standard output (without doing any special redirecting itself, just doing normal `print` statements or whatever), and the next program reads from standard input, and the operating system takes care of connecting one program's output to the next program's input.

**Example 10.12. Chaining commands**

```
[you@localhost kgp]$ python kgp.py -g binary.xml        ❶
01100111
[you@localhost kgp]$ cat binary.xml                     ❷
<?xml version="1.0"?>
```

❸ ❹

❶
❷
❸

❹

❶

❶ This is the `openAnything` function from `toolbox.py`, which you previously examined in
Section 10.1, Abstracting input sources . All you've done is add three lines of code at the beginning
of the function to check if the source is "-"; if so, you return `sys.stdin`. Really, that's it!
Remember, `stdin` is a file–like object with a `read` method, so the rest of the code (in `kgp.py`,
where you call `openAnything`) doesn't change a bit.

## 10.3. Caching node lookups

`kgp.py` 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 of the input documents to build a cache of nodes.

A grammar file defines a series of `ref` elements. Each `ref` contains one or more `p` elements, which can contain a lot
of different things, including `xref`s. Whenever you encounter an `xref`, you look for a corresponding `ref` element
with the same `id` attribute, and choose one of the `ref` element's children and parse it. (You'll see how this random
choice is made in the next section.)

This is how you build up the grammar: define `ref` elements for the smallest pieces, then define `ref` elements which
"include" the first `ref` elements by using `xref`, and so forth. Then you parse the "largest" reference and follow each
`xref`, and eventually output real text. The text you output depends on the (random) decisions you make each time
you fill in an `xref`, so the output is different each time.

This is all very flexible, but there is one downside: performance. When you find an `xref` and need to find the
corresponding `ref` element, you have a problem. The `xref` has an `id` attribute, and you want to find the `ref`
element that has that same `id` attribute, but there is no easy way to do that. The slow way to do it would be to get the
entire list of `ref` elements each time, then manually loop through and look at each `id` attribute. The fast way is to do
that once and build a cache, in the form of a dictionary.

**Example 10.14. `loadGrammar`**

```
def loadGrammar(self, grammar):
```

❶
❷
❸ ❹

❶

❷

❸

❹

```
def do_xref(self, node):
    id = node.attributes["id"].value
    self.parse(self.randomChildElement(self.refs[id]))
```

You'll explore the `randomChildElement` function in the next section.

## 10.4. Finding direct children of a node

Another useful techique when parsing XML documents is finding all the direct child elements of a particular element. For instance, in the grammar files, a `ref` element can have several p

❶ ❷ ❸
❹

❶

❷

❸

❹

### Example 10.17. Class names of parsed XML objects

```
>>> from xml.dom import minidom
>>> xmldoc = minidom.parse('kant.xml')  ❶
>>> xmldoc
<xml.dom.minidom.Document instance at 0x01359DE8>
>>> xmldoc.__class__                     ❷
<class xml.dom.minidom.Document at 0x01105D40>
>>> xmldoc.__class__.__name__            ❸
'Document'
```

❶  Assume for a moment that `kant.xml` is in the current directory.

❷  As you saw in Section 9.2,  Packages , the object returned by parsing an XML document is a `Document` object, as defined in the `minidom.py` in the `xml.dom` package. As you saw in Section 5.4,  Instantiating Classes , `__class__` is built–in attribute of every Python object.

❸  Furthermore, `__name__` 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 you type when you define a class yourself. (See Section 5.3,  Defining Classes .)

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 to your advantage to separate the logic of parsing each node type? The answer is `getattr`, which you first saw in Section 4.4,  Getting Object References With getattr .

### Example 10.18. `parse`, a generic XML node dispatcher

```
def parse(self, node):
    parseMethod = getattr(self, "parse_%s" % node.__class__.__name__)  ❶ ❷
    parseMethod(node)  ❸
```

❶  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're passed a `Document` node, you're constructing the string `'parse_Document'`, and so forth.

❷  Now you can treat that string as a function name, and get a reference to the function itself using `getattr`

❸  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 `parse` dispatcher

```
def parse_Document(self, node):  ❶
    self.parse(node.documentElement)

def parse_Text(self, node):      ❷
    text = node.data
    if self.capitalizeNextWord:
        self.pieces.append(text[0].upper())
        self.pieces.append(text[1:])
        self.capitalizeNextWord = 0
    else:
        self.pieces.append(text)

def parse_Comment(self, node):  ❸
    pass

def parse_Element(self, node):  ❹
    handlerMethod = getattr(self, "do_%s" % node.tagName)
    handlerMethod(node)
```

❶  `parse_Document` is only ever called once, since there is only one `Document` node in an XML document, and only one `Document` object in the parsed XML representation. It simply turns around and parses the root element of the grammar file.

❷  `parse_Text` is called on nodes that represent bits of text. The function itself does some special processing to handle automatic capitalization of the first word of a sentence, but otherwise simply appends the represented text to a list.

❸  `parse_Comment` is just a `pass`, since you don't care about embedded comments in the grammar files. Note, however, that you still need to define the function and explicitly make it do nothing. If the function did not exist, the generic `parse` function would fail as soon as it stumbled on a comment, because it would try to find the non−existent `parse_Comment` function. Defining a separate function for every node type, even ones you don't use, allows the generic `parse` function to stay simple and dumb.

❹  The `parse_Element` method is actually itself a dispatcher, based on the name of the element's tag. The basic idea is the same: take what distinguishes elements from each other (their tag names) and dispatch to a separate function for each of them. You construct a string like `'do_xref'` (for an `<xref>` tag), find a function of that name, and call it. And so forth for each of the other tag names that might be found in the course of parsing a grammar file (`<p>` tags, `<choice>` tags).

In this example, the dispatch functions `parse` and `parse_Element` simply find other methods in the same class. If your processing is very complex (or you have many different tag names), you could break up your code into separate modules, and use dynamic importing to import each module and call whatever functions you needed. Dynamic importing will be discussed in Chapter 16, *Functional Programming*.

❶

❶

❶

❷

```
[you@localhost py]$ python argecho.py --help          ❸
argecho.py
--help
[you@localhost py]$ python argecho.py -m kant.xml ❹
argecho.py
-m
kant.xml
```

❶ 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 later, in Chapter 16, *Functional Programming*. Don't worry about it for now.

❷ Command–line arguments are separated by spaces, and each shows up as a separate element in the `sys.argv` list.

❸ Command–line flags, like `--help`, also show up as their own element in the `sys.argv` list.

❹ To make things even more interesting, some command–line flags themselves take arguments. For instance, here you have a flag (`-m`) which takes an argument (`kant.xml`). Both the flag itself and the flag's argument are simply sequential elements in the `sys.argv` list. No attempt is made to associate one with the other; all you get is a list.

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 easy to actually use it. For simple programs that only take a single argument and have no flags, you can simply use `sys.argv[1]` to access the argument. There's no shame in this; I do it all the time. For more complex programs, you need the `getopt` module.

So what are all those parameters you pass to the `getopt` function? Well, the first one is simply the raw list of command–line flags and arguments (not including the first element, the script name, which you already chopped off before calling the `main` function). The second is the list of short command–line flags that the script accepts.

**"hg:d"**

*−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 debugging). However, the second flag (−g) *must* be followed by an argument, which is the name of the grammar file to read from. In fact it can be a filename or a web address, and you don't know which yet (you'll figure it out later), but you know it has to be *something*. So you tell `getopt` this by putting a colon after the g in that second parameter to the `getopt` function.

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 what the third parameter to `getopt` is for, to specify a list of the long flags that correspond to the short flags you specified in the second parameter.

**["help", "grammar="]**

*−−help*
> print usage summary

*−−grammar ...*
> use specified grammar file or URL

Three things of note here:

1. All long flags are preceded by two dashes on the command line, but you don't include those dashes when calling `getopt`. They are understood.
2. The −−grammar flag must always be followed by an additional argument, just like the −g flag. This is notated by an equals sign, "grammar=".
3. 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; only −d will turn on debugging. But the order of short and long flags needs to be the same, so you'll need to specify all the short flags that *do* have corresponding long flags first, then all the rest of the short flags.

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 `kgp.py`**

```
def main(argv):                                        ❶
    grammar = "kant.xml"
    try:
        opts, args = getopt.getopt(argv, "hg:d", ["help", "grammar="])
    except getopt.GetoptError:
        usage()
        sys.exit(2)
    for opt, arg in opts:                              ❷
```

```
        if opt in ("-h", "--help"):       ❸
            usage()
            sys.exit()
        elif opt == '-d':                  ❹
            global _debug
            _debug = 1
        elif opt in ("-g", "--grammar"):   ❺
            grammar = arg

    source = "".join(args)                 ❻

    k = KantGenerator(grammar, source)
    print k.output()
```

❶   The grammar variable will keep track of the grammar file you're using. You initialize it here in case it's not

❷

❸

❹

❺

❻

```
def _load(self, source):
    sock = toolbox.openAnything(source)
    xmldoc = minidom.parse(sock).documentElement
    sock.close()
```

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 are just elements in the XML document.

```
def loadGrammar(self, grammar):
    for ref in self.grammar.getElementsByTagName("ref"):
        self.refs[ref.attributes["id"].value] = ref
```

If you specified some source material on the command line, you use that; otherwise you rip through the grammar looking for the "top–level" reference (that isn't referenced by anything else) and use that as a starting point.

```
...
    k = KantGenerator(grammar, source)
    print k.output()
```

## 10.8. Summary

Python comes with powerful libraries for parsing and manipulating XML documents. The `minidom` takes an XML file and parses it into Python objects, providing for random access to arbitrary elements. Furthermore, this chapter shows how Python can be used to create a "real" standalone command–line script, complete with command–line

# 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 from a URL, but let's dive into the more general topic of HTTP web services.

Simply stated, HTTP web services are programmatic ways of sending and receiving data from remote servers using the operations of HTTP directly. If you want to get data from the server, use a straight HTTP GET; if you want to send new data to the server, use HTTP POST. (Some more advanced HTTP web service APIs also define ways of modifying existing data and deleting data, using HTTP PUT and HTTP DELETE.) In other words, the "verbs" built into the HTTP protocol (GET, POST, PUT, and DELETE) map directly to application–level operations for receiving, sending, modifying, and deleting data.

The main advantage of this approach is simplicity, and its simplicity has proven popular with a lot of different sites. Data −− usually XML data −− can be built and stored statically, or generated dynamically by a server–side script, and all major languages include an HTTP library for downloading it. Debugging is also easier, because you can load up the web service in any web browser and see the raw data. Modern browsers will even nicely format and pretty–print XML data for you, to allow you to quickly navigate through it.

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/) (Hong Kong) offer weather alerts as a web service.
- 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 underlying HTTP semantics. (They tunnel everything over HTTP POST.) But this chapter will concentrate on using HTTP GET to get data from a remote server, and you'll explore several HTTP features you can use to get the maximum benefit out of pure HTTP web services.

Here is a more advanced version of the `openanything`

```
            return result

    def http_error_302(self, req, fp, code, msg, headers):
        result = urllib2.HTTPRedirectHandler.http_error_302(
            self, req, fp, code, msg, headers)
        result.status = code
        return result

class DefaultErrorHandler(urllib2.HTTPDefaultErrorHandler):
    def http_error_default(self, req, fp, code, msg, headers):
        result = urllib2.HTTPError(
            req.get_full_url(), code, msg, headers, fp)
        result.status = code
        return result

def openAnything(source, etag=None, lastmodified=None, agent=USER_AGENT):
    '''URL, filename, or string --> stream

    This function lets you define parsers that take any input source
    (URL, pathname to local or network file, or actual data as a string)
    and deal with it in a uniform manner.  Returned object is guaranteed
    to have all the basic stdio read methods (read, readline, readlines).
    Just .close() the object when you're done with it.

    If the etag argument is supplied, it will be used as the value of an
    If-None-Match request header.

    If the lastmodified argument is supplied, it must be a formatted
    date/time string in GMT (as returned in the Last-Modified header of
    a previous request).  The formatted date/time will be used
    as the value of an If-Modified-Since request header.

    If the agent argument is supplied, it will be used as the value of a
    User-Agent request header.
    '''

    if hasattr(source, 'read'):
        return source

    if source == '-':
        return sys.stdin

    if urlparse.urlparse(source)[0] == 'http':
        # open URL with urllib2
        request = urllib2.Request(source)
        request.add_header('User-Agent', agent)
        if etag:
            request.add_header('If-None-Match', etag)
        if lastmodified:
            request.add_header('If-Modified-Since', lastmodified)
        request.add_header('Accept-encoding', 'gzip')
        opener = urllib2.build_opener(SmartRedirectHandler(), DefaultErrorHandler())
        return opener.open(request)

    # try to open with native open function (if source is a filename)
    try:
        return open(source)
    except (IOError, OSError):
        pass

    # treat source as string
    return StringIO(str(source))
```

❶

❶

# 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. When the client requests a resource, it should always announce who it is, as specifically as possible. This allows the server–side administrator to get in touch with the client–side developer if anything is going fantastically wrong.

By default, Python sends a generic `User-Agent: Python-urllib/1.15`. 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 at `http://example.com/index.xml` might be moved to `http://example.com/xml/atom.xml`. Or an entire domain might move, as an organization expands and reorganizes; for instance, `http://www.example.com/index.xml` might be redirected to `http://server-farm-1.example.com/index.xml`.

Every time you request any kind of resource from an HTTP server, the server includes a status code in its response. Status code `200` means "everything's normal, here's the page you asked for". Status code `404` means "page not found". (You've probably seen 404 errors while browsing the web.)

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 here temporarily" (and then gives the temporary address in a `Location:` header). Status code `301` is a *permanent redirect*; it means "oops, that got moved permanently" (and then gives the new address in a `Location:` header). If you get a `302` status code and a new address, the HTTP specification says you should use the new address to get what you asked for, but the next time you want to access the same resource, you should retry the old address. But if you get a `301` status code and a new address, you're supposed to use the new address from then on.

`urllib.urlopen` will automatically "follow" redirects when it receives the appropriate status code from the HTTP server, but unfortunately, it doesn't tell you when it does so. You'll end up getting data you asked for, but you'll never know that the underlying library "helpfully" followed a redirect for you. So you'll continue pounding away at the old address, and each time you'll get redirected to the new address. That's two round trips instead of one: not very efficient! Later in this chapter, you'll see how to work around this so you can deal with permanent redirects properly and efficiently.

## 11.3.3. `Last-Modified/If-Modified-Since`

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 Google.com only changes once every few weeks (when they put up a special holiday logo, or advertise a new service). Web services are no different; usually the server knows when the data you requested last changed, and HTTP provides a way for the server to include this last–modified date along with the data you requested.

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 an `If-Modified-Since` header with your request, with the date you got back from the server last time. If the data hasn't changed since then, the server sends back a special HTTP status code `304`, which

means "this data hasn't changed since the last time you asked for it". Why is this an improvement? Because when the server sends a `304`, *it doesn't re−send the data*. All you get is the status code. So you don't need to download the same data over and over again if it hasn't changed; the server assumes you have the data cached locally.

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 hadn't changed, and wondered why it loaded so quickly the second time −− this could be why. Your web browser cached the contents of the page locally the first time, and when you visited the second time, your browser automatically sent the last−modified date it got from the server the first time. The server simply says `304: Not Modified`, so your browser knows to load the page from its cache. Web services can be this smart too.

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 arbitrary headers in each response, you can add support for it yourself.

## 11.3.4. `ETag/If-None-Match`

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 way it works is, the server sends some sort of hash of the data (in an `ETag` header) along with the data you requested. Exactly how this hash is determined is entirely up to the server. The second time you request the same data, you include the ETag hash in an `If-None-Match:` header, and if the data hasn't changed, the server will send you back a `304` status code. As with the last−modified date checking, the server *just* sends the `304`; it doesn't send you the same data a second time. By including the ETag hash in your second request, you're telling the server that there's no need to re−send the same data if it still matches this hash, since you still have the data from the last time.

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 /F9eooi the data fied date checkinows to lznhe l kheckinows st important HTa

```
>>> import urllib
>>> feeddata = urllib.urlopen('http://diveintomark.org/xml/atom.xml').read()
connect: (diveintomark.org, 80)                              ❷
send: '
GET /xml/atom.xml HTTP/1.0                                   ❸
Host: diveintomark.org                                       ❹
User-agent: Python-urllib/1.15                               ❺
'
reply: 'HTTP/1.1 200 OK\r\n'                                 ❻
header: Date: Wed, 14 Apr 2004 22:27:30 GMT
header: Server: Apache/2.0.49 (Debian GNU/Linux)
header: Content-Type: application/atom+xml
header: Last-Modified: Wed, 14 Apr 2004 22:14:38 GMT  ❼
header: ETag: "e8284-68e0-4de30f80"                   ❽
header: Accept-Ranges: bytes
header: Content-Length: 26848
header: Connection: close
```

❶     `urllib` relies on another standard Python library, `httplib`. Normally you don't need to `import httplib` directly (`urllib` does that automatically), but you will here so you can set the debugging flag on the `HTTPConnection` class that `urllib` uses internally to connect to the HTTP server. This is an incredibly useful technique. Some other Python libraries have similar debug flags, but there's no particular standard for naming them or turning them on; you need to read the documentation of each library to see if such a feature is available.

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 that you're connecting to the server `diveintomark.org` on port 80, which is the standard port for HTTP.

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 (minus the domain name). All the requests in this chapter will use `GET`, but in the next chapter on SOAP, you'll see that it uses `POST` for everything. The basic syntax is the same, regardless of the verb.

❹     The second line is the `Host` header, which specifies the domain name of the service you're accessing

❺

❻

❼

❽

## 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 move beyond the basic `urllib` and dive into `urllib2`.

**Example 11.4. Introducing `urllib2`**

```
>>> import httplib
>>> httplib.HTTPConnection.debuglevel = 1                              ❶
>>> import urllib2
>>> request = urllib2.Request('http://diveintomark.org/xml/atom.xml')  ❷
>>> opener = urllib2.build_opener()                                    ❸
>>> feeddata = opener.open(request).read()                            ❹
connect: (diveintomark.org, 80)
send: '
GET /xml/atom.xml HTTP/1.0
Host: diveintomark.org
User-agent: Python-urllib/2.1
'
reply: 'HTTP/1.1 200 OK\r\n'
header: Date: Wed, 14 Apr 2004 23:23:12 GMT
header: Server: Apache/2.0.49 (Debian GNU/Linux)
header: Content-Type: application/atom+xml
header: Last-Modified: Wed, 14 Apr 2004 22:14:38 GMT
header: ETag: "e8284-68e0-4de30f80"
header: Accept-Ranges: bytes
header: Content-Length: 26848
header: Connection: close
```

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 you're actually sending over the wire, and what gets sent back.

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 `Request` object, which takes the URL of the resource you'll eventually get around to retrieving. Note that this step doesn't actually retrieve anything yet.

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 build an opener without any custom handlers, which is what you're doing here. You'll see how to define and use custom handlers later in this chapter when you explore redirects.

❹   The final step is to tell the opener to open the URL, using the `Request` object you created. As you can see from all the debugging information that gets printed, this step actually retrieves the resource and stores the returned data in `feeddata`.

**Example 11.5. Adding headers with the `Request`**

```
>>> request                                              ❶
<urllib2.Request instance at 0x00250AA8>
>>> request.get_full_url()
http://diveintomark.org/xml/atom.xml
>>> request.add_header('User-Agent',
...      'OpenAnything/1.0 +http://diveintopython.org/')  ❷
>>> feeddata = opener.open(request).read()                ❸
connect: (diveintomark.org, 80)
send: '
GET /xml/atom.xml HTTP/1.0
Host: diveintomark.org
User-agent: OpenAnything/1.0 +http://diveintopython.org/  ❹
'
```

```
reply: 'HTTP/1.1 200 OK\r\n'
header: Date: Wed, 14 Apr 2004 23:45:17 GMT
header: Server: Apache/2.0.49 (Debian GNU/Linux)
header: Content-Type: application/atom+xml
header: Last-Modified: Wed, 14 Apr 2004 22:14:38 GMT
header: ETag: "e8284-68e0-4de30f80"
header: Accept-Ranges: bytes
header: Content-Length: 26848
header: Connection: close
```

❶     You're continuing from the previous example; you've already created a `Request` object with the URL you want to access.

❷     Using the `add_header` method on the `Request` object, you can add arbitrary HTTP headers to the

❸

❹

❶

❷
❸

```
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
  File "c:\python23\lib\urllib2.py", line 326, in open
    '_open', req)
  File "c:\python23\lib\urllib2.py", line 306, in _call_chain
    result = func(*args)
  File "c:\python23\lib\urllib2.py", line 901, in http_open
    return self.do_open(httplib.HTTP, req)
  File "c:\python23\lib\urllib2.py", line 895, in do_open
    return self.parent.error('http', req, fp, code, msg, hdrs)
  File "c:\python23\lib\urllib2.py", line 352, in error
    return self._call_chain(*args)
  File "c:\python23\lib\urllib2.py", line 306, in _call_chain
    result = func(*args)
  File "c:\python23\lib\urllib2.py", line 412, in http_error_default
    raise HTTPError(req.get_full_url(), code, msg, hdrs, fp)
urllib2.HTTPError: HTTP Error 304: Not Modified
```

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: `firstdatastream.headers` is an object that acts like a dictionary and allows you to get any of the individual headers returned from the HTTP server.

❷　On the second request, you add the `If-Modified-Since` header with the last–modified date from the first

❸

❶
❷

❸

❶

❷

❸

**Example 11.8. Using custom URL handlers**

```
>>> request.headers                              ❶
{'If-modified-since': 'Thu, 15 Apr 2004 19:45:21 GMT'}
>>> import openanything
>>> opener = urllib2.build_opener(
...     openanything.DefaultErrorHandler())     ❷
>>> seconddatastream = opener.open(request)
>>> seconddatastream.status                      ❸
304
>>> seconddatastream.read()                      ❹
''
```

❶   You're continuing the previous example, so the `Request` object is already set up, and you've already added the `If-Modified-Since` header.

❷   This is the key: now that you've defined your custom URL handler, you need to tell `urllib2` to use it. Remember how I said that `urllib2` broke up the process of accessing an HTTP resource into three steps, and for good reason? This is why building the URL opener is its own step, because you can build it with your own custom URL handlers that override `urllib2`'s default behavior.

❸   Now you can quietly open the resource, and what you get back is an object that, along with the usual headers (use `seconddatastream.headers.dict` to acess them), also contains the HTTP status code. In this case, as you expected, the status is `304`, meaning this data hasn't changed since the last time you asked for it.

❹   Note that when the server sends back a `304` status code, it doesn't re–send the data. That's the whole point: to save bandwidth by not re–downloading data that hasn't changed. So if you actually want that data, you'll need to cache it locally the first time you get it.

Handling `ETag` works much the same way, but instead of checking for `Last-Modified` and sending `If-Modified-Since`, you check for `ETag` and send `If-None-Match`. Let's start with a fresh IDE session.

**Example 11.9. Supporting `ETag/If-None-Match`**

```
>>> import urllib2, openanything
>>> request = urllib2.Request('http://diveintomark.org/xml/atom.xml')
>>> opener = urllib2.build_opener(
...     openanything.DefaultErrorHandler())
>>> firstdatastream = opener.open(request)
>>> firstdatastream.headers.get('ETag')          ❶
'"e842a-3e53-55d97640"'
>>> firstdata = firstdatastream.read()
>>> print firstdata                              ❷
<?xml version="1.0" encoding="iso-8859-1"?>
<feed version="0.3"
  xmlns="http://purl.org/atom/ns#"
  xmlns:dc="http://purl.org/dc/elements/1.1/"
  xml:lang="en">
  <title mode="escaped">dive into mark</title>
  <link rel="alternate" type="text/html" href="http://diveintomark.org/"/>
  <-- rest of feed omitted for brevity -->
>>> request.add_header('If-None-Match',
...     firstdatastream.headers.get('ETag'))     ❸
>>> seconddatastream = opener.open(request)
>>> seconddatastream.status                       ❹
304
>>> seconddatastream.read()                       ❺
''
```

❶

Using the `firstdatastream.headers` pseudo–dictionary, you can get the `ETag`
returned from the server. (What happens if the server didn't send back an `ETag`? Then this line
would return `None`.)

❷ OK, you got the data.

❸ Now set up the second call by setting the `If-None-Match` header to the `ETag` you got from
the first call.

❹ The second call succeeds quietly (without throwing an exception), and once again you see that
the server has sent back a `304` status code. Based on the `ETag` you sent the second time, it
knows that the data hasn't changed.

❺ Regardless of whether the `304` is triggered by `Last-Modified` date checking or `ETag`
hash matching, you'll never get the data along with the `304`. That's the whole point.

In these examples, the HTTP server has supported both `Last-Modified` and `ETag` headers, but not all servers do.
As a web services client, you should be prepared to support both, but you must code defensively in case a server only
supports one or the other, or neither.

# 11.7. Handling redirects

❶

❷

❸

❹

❺

```
header: Content-Type: application/atom+xml
>>> f.url                                                    ❻
'http://diveintomark.org/xml/atom.xml'
>>> f.headers.dict
{'content-length': '15955',
'accept-ranges': 'bytes',
'server': 'Apache/2.0.49 (Debian GNU/Linux)',
'last-modified': 'Thu, 15 Apr 2004 19:45:21 GMT',
'connection': 'close',
'etag': '"e842a-3e53-55d97640"',
'date': 'Thu, 15 Apr 2004 22:06:25 GMT',
'content-type': 'application/atom+xml'}
>>> f.status
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
AttributeError: addinfourl instance has no attribute 'status'
```

❶   You'll be better able to see what's happening if you turn on debugging.

❷   This is a URL which I have set up to permanently redirect to my Atom feed at
    `http://diveintomark.org/xml/atom.xml`.

❸   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 permanently.

❹   The server also sends back a `Location:` header that gives the new address of this data.

❺   `urllib2` notices the redirect status code and automatically tries to retrieve the data at the new location
    specified in the `Location:` header.

❻   The object you get back from the `opener` contains the new permanent address and all the headers returned
    from the second request (retrieved from the new permanent address). But the status code is missing, so you
    have no way of knowing programmatically whether this redirect was temporary or permanent. And that matters
    very much: if it was a temporary redirect, then you should continue to ask for the data at the old location. But if
    it was a permanent redirect (as this was), you should ask for the data at the new location from now on.

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? With a custom URL handler, just like you did to handle `304` codes.


**Example 11.11. Defining the redirect handler**

This class is defined in `openanything.py`.

```
class SmartRedirectHandler(urllib2.HTTPRedirectHandler):       ❶
    def http_error_301(self, req, fp, code, msg, headers):
        result = urllib2.HTTPRedirectHandler.http_error_301(   ❷
            self, req, fp, code, msg, headers)
        result.status = code                                   ❸
        return result

    def http_error_302(self, req, fp, code, msg, headers):     ❹
        result = urllib2.HTTPRedirectHandler.http_error_302(
            self, req, fp, code, msg, headers)
        result.status = code
        return result
```

❶   Redirect behavior is defined in `urllib2` in a class called `HTTPRedirectHandler`. You
    don't want to completely override the behavior, you just want to extend it a little, so you'll
    subclass `HTTPRedirectHandler` so you can call the ancestor class to do all the hard work.

❷

When it encounters a `301` status code from the server, `urllib2` will search through its handlers and call the `http_error_301` method. The first thing ours does is just call the `http_error_301` method in the ancestor, which handles the grunt work of looking for the `Location:` header and following the redirect to the new address.

❸  Here's the key: before you return, you store the status code (`301`), so that the calling program can access it later.

❹  Temporary redirects (status code `302`) work the same way: override the `http_error_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, but now it will also expose the redirect status code.

**Example 11.12. Using the redirect handler to detect permanent redirects**

```
>>> request = urllib2.Request('http://diveintomark.org/redir/example301.xml')
>>> import openanything, httplib
>>> httplib.HTTPConnection.debuglevel = 1
>>> opener = urllib2.build_opener(
...     openanything.SmartRedirectHandler())          ❶
>>> f = opener.open(request)
connect: (diveintomark.org, 80)
send: 'GET /redir/example301.xml HTTP/1.0
Host: diveintomark.org
User-agent: Python-urllib/2.1
'
reply: 'HTTP/1.1 301 Moved Permanently\r\n'          ❷
header: Date: Thu, 15 Apr 2004 22:13:21 GMT
header: Server: Apache/2.0.49 (Debian GNU/Linux)
header: Location: http://diveintomark.org/xml/atom.xml
header: Content-Length: 338
header: Connection: close
header: Content-Type: text/html; charset=iso-8859-1
connect: (diveintomark.org, 80)
send: '
GET /xml/atom.xml HTTP/1.0
Host: diveintomark.org
User-agent: Python-urllib/2.1
'
reply: 'HTTP/1.1 200 OK\r\n'
header: Date: Thu, 15 Apr 2004 22:13:21 GMT
header: Server: Apache/2.0.49 (Debian GNU/Linux)
header: Last-Modified: Thu, 15 Apr 2004 19:45:21 GMT
header: ETag: "e842a-3e53-55d97640"
header: Accept-Ranges: bytes
header: Content-Length: 15955
header: Connection: close
header: Content-Type: application/atom+xml

>>> f.status                                          ❸
301
>>> f.url
'http://diveintomark.org/xml/atom.xml'
```

❶  First, build a URL opener with the redirect handler you just defined.

❷  You sent off a request, and you got a `301` status code in response. At this point, the `http_error_301` method gets called. You call the ancestor method, which follows the redirect and sends a request at the new location (`http://diveintomark.org/xml/atom.xml`).

❸     This is the payoff: now, not only do you have access to the new URL, but you have access to the redirect status

❶

❷

❸

❹

❶

❷

❸

❹

maybe not. Maybe it will redirect to a different address. It's not for you to say. The server said this redirect was only temporary, so you should respect that. And now you're exposing enough information that the calling application can respect that.

## 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 cut down the amount of data sent over the wire by 60% or more. This is especially true of XML web services, since XML data compresses very well.

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**

```
>>> import urllib2, httplib
>>> httplib.HTTPConnection.debuglevel = 1
>>> request = urllib2.Request('http://diveintomark.org/xml/atom.xml')
>>> request.add_header('Accept-encoding', 'gzip')          ❶
>>> opener = urllib2.build_opener()
>>> f = opener.open(request)
connect: (diveintomark.org, 80)
send: '
GET /xml/atom.xml HTTP/1.0
Host: diveintomark.org
User-agent: Python-urllib/2.1
Accept-encoding: gzip                                      ❷
'
reply: 'HTTP/1.1 200 OK\r\n'
header: Date: Thu, 15 Apr 2004 22:24:39 GMT
header: Server: Apache/2.0.49 (Debian GNU/Linux)
header: Last-Modified: Thu, 15 Apr 2004 19:45:21 GMT
header: ETag: "e842a-3e53-55d97640"
header: Accept-Ranges: bytes
header: Vary: Accept-Encoding
header: Content-Encoding: gzip                             ❸
header: Content-Length: 6289                               ❹
header: Connection: close
header: Content-Type: application/atom+xml
```

❶    This is the key: once you've created your `Request` object, add an `Accept-encoding` header to tell the server you can accept gzip–encoded data. `gzip` is the name of the compression algorithm you're using. In theory there could be other compression algorithms, but `gzip` is the compression algorithm used by 99% of web servers.

❷    There's your header going across the wire.

❸    And here's what the server sends back: the `Content-Encoding: gzip` header means that the data you're about to receive has been gzip–compressed.

❹    The `Content-Length` 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 was 15955, so gzip compression cut your bandwidth by over 60%!

**Example 11.15. Decompressing the data**

```
>>> compresseddata = f.read()                              ❶
>>> len(compresseddata)
6289
>>> import StringIO
```

```
>>> compressedstream = StringIO.StringIO(compresseddata)      ❷
>>> import gzip
>>> gzipper = gzip.GzipFile(fileobj=compressedstream)          ❸
>>> data = gzipper.read()                                      ❹
>>> print data                                                 ❺
<?xml version="1.0" encoding="iso-8859-1"?>
<feed version="0.3"
  xmlns="http://purl.org/atom/ns#"
  xmlns:dc="http://purl.org/dc/elements/1.1/"
  xml:lang="en">
  <title mode="escaped">dive into mark</title>
  <link rel="alternate" type="text/html" href="http://diveintomark.org/"/>
  <-- rest of feed omitted for brevity -->
>>> len(data)
15955
```

❶   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 uncompressed data, but since this data has been gzip–compressed, this is just the first step towards getting the data you really want.

❷   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. But you don't have a file on disk, you have a gzip–compressed buffer in memory, and you don't want to write out a temporary file just so you can uncompress it. So what you're going to do is create a file–like object out of the in–memory data (`compresseddata`), using the `StringIO` module. You first saw the `StringIO` module in the previous chapter, but now you've found another use for it.

❸   Now you can create an instance of `GzipFile`, and tell it that its "file" is the file–like object `compressedstream`.

❹   This is the line that does all the actual work: "reading" from `GzipFile` will decompress the data. Strange? Yes, but it makes sense in a twisted kind of way. `gzipper` is a file–like object which represents a gzip–compressed file. That "file" is not a real file on disk, though; `gzipper` is really just "reading" from the file–like object you created with `StringIO` to wrap the compressed data, which is only in memory in the variable `compresseddata`. And where did that compressed data come from? You originally downloaded it from a remote HTTP server by "reading" from the file–like object you built with `urllib2.build_opener`. And amazingly, this all just works. Every step in the chain has no idea that the previous step is faking it.

❺   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 `opener.open` returns a file–like object, so why not cut out the `StringIO` middleman and just pass `f` directly to `GzipFile`? OK, maybe you weren't thinking that, but don't worry about it, because it doesn't work.

**Example 11.16. Decompressing the data directly from the server**

```
>>> f = opener.open(request)                          ❶
>>> f.headers.get('Content-Encoding')                 ❷
'gzip'
>>> data = gzip.GzipFile(fileobj=f).read()            ❸
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
  File "c:\python23\lib\gzip.py", line 217, in read
    self._read(readsize)
  File "c:\python23\lib\gzip.py", line 252, in _read
    pos = self.fileobj.tell()   # Save current position
AttributeError: addinfourl instance has no attribute 'tell'
```

❶ Continuing from the previous example, you already have a `Request` object set up with an `Accept-encoding: gzip` header.

❷ Simply opening the request will get you the headers (though not download any data yet). As you can see from the returned `Content-Encoding` header, this data has been sent gzip–compressed.

❸ Since `opener.open`

❶

❷

❸

❹
❺
❻
❼

❶

❷

❸

❹

❺

❻

❼

This function is defined in `openanything.py`.

```
def fetch(source, etag=None, last_modified=None, agent=USER_AGENT):
    '''Fetch data and metadata from a URL, file, stream, or string'''
    result = {}
    f = openAnything(source, etag, last_modified, agent)          ❶
    result['data'] = f.read()                                     ❷
    if hasattr(f, 'headers'):
        # save ETag, if the server sent one
        result['etag'] = f.headers.get('ETag')                    ❸
        # save Last-Modified header, if the server sent one
        result['lastmodified'] = f.headers.get('Last-Modified')   ❹
        if f.headers.get('content-encoding', '') == 'gzip':       ❺
            # data came back gzip-compressed, decompress it
            result['data'] = gzip.GzipFile(fileobj=StringIO(result['data'])).read()
    if hasattr(f, 'url'):                                          ❻
        result['url'] = f.url
        result['status'] = 200
    if hasattr(f, 'status'):                                       ❼
        result['status'] = f.status
    f.close()
    return result
```

❶   First, you call the `openAnything` function with a URL, `ETag` hash, `Last-Modified` date, and `User-Agent`.

❷   Read the actual data returned from the server. This may be compressed; if so, you'll decompress it later.

❸   Save the `ETag` hash returned from the server, so the calling application can pass it back to you next time, and you can pass it on to `openAnything`, which can stick it in the `If-None-Match` header and send it to the remote server.

❹   Save the `Last-Modified` date too.

❺   If the server says that it sent compressed data, decompress it.

❻   If you got a URL back from the server, save it, and assume that the status code is `200` until you find out otherwise.

❼   If one of the custom URL handlers captured a status code, then save that too.

**Example 11.19. Using `openanything.py`**

```
>>> import openanything
>>> useragent = 'MyHTTPWebServicesApp/1.0'
>>> url = 'http://diveintopython.org/redir/example301.xml'
>>> params = openanything.fetch(url, agent=useragent)          ❶
>>> params                                                     ❷
{'url': 'http://diveintomark.org/xml/atom.xml',
'lastmodified': 'Thu, 15 Apr 2004 19:45:21 GMT',
'etag': '"e842a-3e53-55d97640"',
'status': 301,
'data': '<?xml version="1.0" encoding="iso-8859-1"?>
<feed version="0.3"
<-- rest of data omitted for brevity -->'}
>>> if params['status'] == 301:                                ❸
...     url = params['url']
>>> newparams = openanything.fetch(
...     url, params['etag'], params['lastmodified'], useragent) ❹
>>> newparams
{'url': 'http://diveintomark.org/xml/atom.xml',
'lastmodified': None,
'etag': '"e842a-3e53-55d97640"',
'status': 304,
```

`'data': ''}`  ❺

❶    The very first time you fetch a resource, you don't have an `ETag` hash or `Last-Modified` date, 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. `openanything` handles the gzip compression internally; you don't care about that at this level.

❸    If you ever get a `301`

❹

❺

# 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 XML document which it was your responsibility to parse.

This chapter will focus on SOAP web services, which take a more structured approach. Rather than dealing with HTTP requests and XML documents directly, SOAP allows you to simulate calling functions that return native data types. As you will see, the illusion is almost perfect; you can "call" a function through a SOAP library, with the standard Python calling syntax, and the function appears to return Python objects and values. But under the covers, the SOAP library has actually performed a complex transaction involving multiple XML documents and a remote server.

SOAP is a complex specification, and it is somewhat misleading to say that SOAP is all about calling remote functions. Some people would pipe up to add that SOAP allows for one–way asynchronous message passing, and document–oriented web services. And those people would be correct; SOAP can be used that way, and in many different ways. But this chapter will focus on so–called "RPC–style" SOAP –– calling a remote function and getting results back.

## 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. Here is a program to search Google from Python.

**Example 12.1. `search.py`**

```
from SOAPpy import WSDL

# you'll need to configure these two values;
# see http://www.google.com/apis/
WSDLFILE = '/path/to/copy/of/GoogleSearch.wsdl'
APIKEY = 'YOUR_GOOGLE_API_KEY'

_server = WSDL.Proxy(WSDLFILE)
def search(q):
    """Search Google and return list of {title, link, description}"""
    results = _server.doGoogleSearch(
        APIKEY, q, 0, 10, False, "", False, "", "utf-8", "utf-8")
    return [{"title": r.title.encode("utf-8"),
             "link": r.URL.encode("utf-8"),
             "description": r.snippet.encode("utf-8")}
            for r in results.resultElements]

if __name__ == '__main__':
    import sys
    for r in search(sys.argv[1])[:5]:
        print r['title']
        print r['link']
        print r['description']
        print
```

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 give the search query as a command–line argument, and it prints out the URL, title, and description of the top five Google search results.

Here is the sample output for a search for the word "python".

will be `PyXML-0.8.3.win32-py2.3.exe`.

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 Start Menu or shortcuts installed on the desktop). PyXML is simply a collection of XML libraries used by other programs.

To verify that you installed PyXML correctly, run your Python IDE and check the version of the XML libraries you have installed, as shown here.

**Procedure 12.3.**

Here is the procedure for installing SOAPpy:

1. 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 `.zip` file; otherwise, download the `.tar.gz` 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

```
>>> import SOAPpy
>>> SOAPpy.__version__
'0.11.4'
```

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 demonstration purposes.

The most popular public access SOAP server is http://www.xmethods.net/. This example uses a demonstration function that takes a United States zip code and returns the current temperature in that region.

### Example 12.6. Getting the Current Temperature

```
>>> from SOAPpy import SOAPProxy              ❶
>>> url = 'http://services.xmethods.net:80/soap/servlet/rpcrouter'
>>> namespace = 'urn:xmethods-Temperature'   ❷
>>> server = SOAPProxy(url, namespace)        ❸
>>> server.getTemp('27502')                   ❹
80.0
```

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 XML request document out of the function name and argument list, sending the request over HTTP to the remote SOAP server, parsing the XML response document, and creating native Python values to return. You'll see what these XML documents look like in the next section.

❷ Every SOAP service has a URL which handles all the requests. The same URL is used for all function calls. This particular service only has a single function, but later in this chapter you'll see examples of the Google API, which has several functions. The service URL is shared by all functions.Each SOAP service also has a namespace, which is defined by the server and is completely arbitrary. It's simply part of the configuration required to call SOAP methods. It allows the server to share a single service URL and route requests between several unrelated services. It's like dividing Python modules into packages.

❸ You're creating the `SOAPProxy` with the service URL and the service namespace. This doesn't make any connection to the SOAP server; it simply creates a local Python object.

❹　Now with everything configured properly, you can actually call remote SOAP methods as if they were local functions. You pass arguments just like a normal function, and you get a return value just like a normal function. But under the covers, there's a heck of a lot going on.

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 `SOAPProxy`'s configuration.

**Example 12.7. Debugging SOAP Web Services**

```
>>> from SOAPpy import SOAPProxy
>>> url = 'http://services.xmethods.net:80/soap/servlet/rpcrouter'
>>> n = 'urn:xmethods-Temperature'
>>> server = SOAPProxy(url, namespace=n)          ❶
>>> server.config.dumpSOAPOut = 1'27502'(serverDo 6 D5 6 D5 6rgek u (from*** 63 es prorvice**********
>"> serve<ike a nxsi:type="xsd:float">8 Do</ike a >>> serve</ns1:jus 0 -R>>> the>>> serve</rvic-ENV:
                                         ❸
```

❶
❷

❸

document, and the incoming XML response document. This is all the hard work that `SOAPProxy` is doing for you. Intimidating, isn't it? Let's break it down.

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 the same (or similar) for all SOAP calls. The heart of the "function call" is this fragment within the `<Body>` element:

```
<ns1:getTemp                                    ❶
  xmlns:ns1="urn:xmethods-Temperature"          ❷
  SOAP-ENC:root="1">
<v1 xsi:type="xsd:string">27502</v1>            ❸
</ns1:getTemp>
```

❶ The element name is the function name, `getTemp`. `SOAPProxy` uses `getattr` as a dispatcher. Instead of calling separate local methods based on the method name, it actually uses the method name to construct the XML request document.

❷ The function's XML element is contained in a specific namespace, which is the namespace you specified when you created the `SOAPProxy` object. Don't worry about the `SOAP-ENC:root`; that's boilerplate too.

❸ The arguments of the function also got translated into XML. `SOAPProxy` introspects each argument to determine its datatype (in this case it's a string). The argument datatype goes into the `xsi:type` attribute, followed by the actual string value.

The XML return document is equally easy to understand, once you know what to ignore. Focus on this fragment within the `<Body>`:

```
<ns1:getTempResponse                                        ❶
  xmlns:ns1="urn:xmethods-Temperature"                      ❷
  SOAP-ENV:encodingStyle="http://schemas.xmlsoap.org/soap/encoding/">
<return xsi:type="xsd:float">80.0</return>                  ❸
</ns1:getTempResponse>
```

❶ The server wraps the function return value within a `<getTempResponse>` element. By convention, this deepermelemerittis methrains of theyfunctionc plusnResponsefolwhich is wu s –13.2Tj lgeou created the

❷

❸

The big difference is introspection. As you saw in Chapter 4, Python excels at letting you discover things about modules and functions at runtime. You can list the available functions within a module, and with a little work, drill down to individual function declarations and arguments.

WSDL lets you do that with SOAP web services. WSDL stands for "Web Services Description Language". Although designed to be flexible enough to describe many types of web services, it is most often used to describe SOAP web services.

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 services it describes, although there's nothing special about it. Later in this chapter, we'll download the WSDL file for the Google API and use it locally. That doesn't mean we're calling Google locally; the WSDL file still describes the remote functions sitting on Google's server.

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 variety of web services)
- 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 history entirely, since it bores me to tears. There were other standards that tried to do similar things, but WSDL won, so let's learn how to use it.

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**

```
>>> from SOAPpy import WSDL              ❶
>>> wsdlFile = 'http://www.xmethods.net/sd/2001/TemperatureService.wsdl')
>>> server = WSDL.Proxy(wsdlFile)        ❷
>>> server.methods.keys()                ❸
[u'getTemp']
```

❶  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 parsing any of the WSDL files I tried.

❷  To use a WSDL file, you again use a proxy class, `WSDL.Proxy`, which takes a single argument: the WSDL file. Note that in this case you are passing in the URL of a WSDL file stored on the remote server, but the proxy class works just as well with a local copy of the WSDL file. The act of creating the WSDL proxy will download the WSDL file and parse it, so it there are any errors in the WSDL file (or it can't be fetched due to networking problems), you'll know about it immediately.

❸  The WSDL proxy class exposes the available functions as a Python dictionary, `server.methods`. So getting the list of available methods is as simple as calling the dictionary method `keys()`.

Okay, so you know that this SOAP server offers a single method: `getTemp`. But how do you call it? The WSDL proxy object can tell you that too.

**Example 12.9. Discovering A Method's Arguments**

```
>>> callInfo = server.methods['getTemp']          ❶
>>> callInfo.inparams                              ❷
[<SOAPpy.wstools.WSDLTools.ParameterInfo instance at 0x00CF3AD0>]
>>> callInfo.inparams[0].name                      ❸
u'zipcode'
>>> callInfo.inparams[0].type                      ❹
(u'http://www.w3.org/2001/XMLSchema', u'string')
```

❶ The `server.methods` dictionary is filled with a SOAPpy–specific structure called `CallInfo`. A `CallInfo` object contains information about one specific function, including the function arguments.

❷ The function arguments are stored in `callInfo.inparams`, which is a Python list of `ParameterInfo` objects that hold information about each parameter.

❸ Each `ParameterInfo` object contains a `name` attribute, which is the argument name. You are not required to know the argument name to call the function through SOAP, but SOAP does support calling functions with named arguments (just like Python), and `WSDL.Proxy` will correctly handle mapping named arguments to the remote function if you choose to use them.

❹ Each parameter is also explicitly typed, using datatypes defined in XML Schema. You saw this in the wire trace in the previous section; the XML Schema namespace was part of the "boilerplate" I told you to ignore. For our purposes here, you may continue to ignore it. The `zipcode` parameter is a string, and if you pass in a Python string to the `WSDL.Proxy` object, it will map it correctly and send it to the server.

WSDL also lets you introspect into a function's return values.

**Example 12.10. Discovering A Method's Return Values**

```
>>> callInfo.outparams                             ❶
[<SOAPpy.wstools.WSDLTools.ParameterInfo instance at 0x00CF3AF8>]
>>> callInfo.outparams[0].name                     ❷
u'return'
>>> callInfo.outparams[0].type
(u'http://www.w3.org/2001/XMLSchema', u'float')
```

❶ The adjunct to `callInfo.inparams` for function arguments is `callInfo.outparams` for return value. It is also a list, because functions called through SOAP can return multiple values, just like Python functions.

❷ Each `ParameterInfo` object contains `name` and `type`. This function returns a single value, named `return`, 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**

```
>>> from SOAPpy import WSDL
>>> wsdlFile = 'http://www.xmethods.net/sd/2001/TemperatureService.wsdl')
>>> server = WSDL.Proxy(wsdlFile)                  ❶
>>> server.getTemp('90210')                        ❷
66.0
>>> server.soapproxy.config.dumpSOAPOut = 1        ❸
>>> server.soapproxy.config.dumpSOAPIn = 1
>>> temperature = server.getTemp('90210')
```

```
*** Outgoing SOAP ****************************************************
<?xml version="1.0" encoding="UTF-8"?>
<SOAP-ENV:Envelope SOAP-ENV:encodingStyle="http://schemas.xmlsoap.org/soap/encoding/"
  xmlns:SOAP-ENC="http://schemas.xmlsoap.org/soap/encoding/"
  xmlns:xsi="http://www.w3.org/1999/XMLSchema-instance"
  xmlns:SOAP-ENV="http://schemas.xmlsoap.org/soap/envelope/"
  xmlns:xsd="http://www.w3.org/1999/XMLSchema">
<SOAP-ENV:Body>
<ns1:getTemp xmlns:ns1="urn:xmethods-Temperature" SOAP-ENC:root="1">
<v1 xsi:type="xsd:string">90210</v1>
</ns1:getTemp>
</SOAP-ENV:Body>
</SOAP-ENV:Envelope>
**********************************************************************
*** Incoming SOAP ****************************************************
<?xml version='1.0' encoding='UTF-8'?>
<SOAP-ENV:Envelope xmlns:SOAP-ENV="http://schemas.xmlsoap.org/soap/envelope/"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xmlns:xsd="http://www.w3.org/2001/XMLSchema">
<SOAP-ENV:Body>
<ns1:getTempResponse xmlns:ns1="urn:xmethods-Temperature"
  SOAP-ENV:encodingStyle="http://schemas.xmlsoap.org/soap/encoding/">
<return xsi:type="xsd:float">66.0</return>
</ns1:getTempResponse>

</SOAP-ENV:Body>
</SOAP-ENV:Envelope>
**********************************************************************

>>> temperature
66.0
```

❶ The configuration is simpler than calling the SOAP service directly, since the WSDL file contains the both service URL and namespace you need to call the service. Creating the `WSDL.Proxy` object downloads the WSDL file, parses it, and configures a `SOAPProxy` object that it uses to call the actual SOAP web service.

❷ Once the `WSDL.Proxy` object is created, you can call a function as easily as you did with the `SOAPProxy` object. This is not surprising; the `WSDL.Proxy` is just a wrapper around the `SOAPProxy` with some introspection methods added, so the syntax for calling functions is the same.

❸ You can access the `WSDL.Proxy`'s `SOAPProxy` with `server.soapproxy`. This is useful to turning on debugging, so that when you can call functions through the WSDL proxy, its `SOAPProxy` will dump the outgoing and incoming XML documents that are going over the wire.

# 12.7. Searching Google

Let's finally turn to the sample code that you saw that the beginning of this chapter, which does something more useful and exciting than get the current temperature.

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 your Google API license key by email. You will need this key to pass as a parameter whenever you call Google's search functions.
2. Also on http://www.google.com/apis/, download the Google Web APIs developer kit. This includes some sample code in several programming languages (but not Python), and more importantly, it includes the WSDL

file.

3. Decompress the developer kit file and find `GoogleSearch.wsdl`. 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**

```
>>> from SOAPpy import WSDL
>>> server = WSDL.Proxy('/path/to/your/GoogleSearch.wsdl')  ❶
>>> server.methods.keys()                                    ❷
[u'doGoogleSearch', u'doGetCachedPage', u'doSpellingSuggestion']
>>> callInfo = server.methods['doGoogleSearch']
>>> for arg in callInfo.inparams:                            ❸
...     print arg.name.ljust(15), arg.type
key             (u'http://www.w3.org/2001/XMLSchema', u'string')
q               (u'http://www.w3.org/2001/XMLSchema', u'string')
start           (u'http://www.w3.org/2001/XMLSchema', u'int')
maxResults      (u'http://www.w3.org/2001/XMLSchema', u'int')
filter          (u'http://www.w3.org/2001/XMLSchema', u'boolean')
restrict        (u'http://www.w3.org/2001/XMLSchema', u'string')
safeSearch      (u'http://www.w3.org/2001/XMLSchema', u'boolean')
lr              (u'http://www.w3.org/2001/XMLSchema', u'string')
ie              (u'http://www.w3.org/2001/XMLSchema', u'string')
oe              (u'http://www.w3.org/2001/XMLSchema', u'string')
```

❶ Getting started with Google web services is easy: just create a `WSDL.Proxy` object and point it at your local copy of Google's WSDL file.

❷ According to the WSDL file, Google offers three functions: `doGoogleSearch`, `doGetCachedPage`, and `doSpellingSuggestion`. These do exactly what they sound like: perform a Google search and return the results programmatically, get access to the cached version of a page from the last time Google saw it, and offer spelling suggestions for commonly misspelled search words.

❸ The `doGoogleSearch` function takes a number of parameters of various types. Note that while the WSDL file can tell you what the arguments are called and what datatype they are, it can't tell you what they mean or how to use them. It could theoretically tell you the acceptable range of values for each parameter, if only specific values were allowed, but Google's WSDL file is not that detailed. `WSDL.Proxy` can't work magic; it can only give you the information provided in the WSDL file.

Here is a brief synopsis of all the parameters to the `doGoogleSearch` function:

- `key` – 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 syntax or tricks, they all work here as well.
- `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 wanted to get the second "page" of results, you would set `start` to 10.
- `maxResults` – The number of results to return. Currently capped at 10, although you can specify fewer if you are only interested in a few results and want to save a little bandwidth.
- `filter` – If `True`, Google will filter out duplicate pages from the results.
  `restrict` – Set this to `country` C7uea=j /F0 11 Tf 7.81 frtrox.Nt0s't

❶

construct the link to the directory category page.

## 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, and receiving XML documents from, an HTTP server. So all the debugging techniques you learned in Chapter 11, *HTTP Web Services* come into play here. You can **import httplib** and then set

❶
❷

❶

❷

```
>>> server = WSDL.Proxy(wsdlFile)
>>> temperature = server.getTemp(27502)                          ❶
<Fault SOAP-ENV:Server: Exception while handling service request:
services.temperature.TempService.getTemp(int) -- no signature match>   ❷
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
  File "c:\python23\Lib\site-packages\SOAPpy\Client.py", line 453, in __call__
    return self.__r_call(*args, **kw)
  File "c:\python23\Lib\site-packages\SOAPpy\Client.py", line 475, in __r_call
    self.__hd, self.__ma)
  File "c:\python23\Lib\site-packages\SOAPpy\Client.py", line 389, in __call
    raise p
SOAPpy.Types.faultType: <Fault SOAP-ENV:Server: Exception while handling service request:
services.temperature.TempService.getTemp(int) -- no signature match>
```

❶  Did you spot the mistake? It's a subtle one: you're calling `server.getTemp` with an integer instead of a string. As you saw from introspecting the WSDL file, the `getTemp()` SOAP function takes a single argument, `zipcode`, which must be a string. `WSDL.Proxy` will *not* coerce datatypes for you; you need to pass the exact datatypes that the server expects.

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 `getTemp` function with an integer value, but there is no function defined with that name that takes an integer. In theory, SOAP allows you to *overload* functions, so you could have two functions in the same SOAP service with the same name and the same number of arguments, but the arguments were of different datatypes. This is why it's important to match the datatypes exactly, and why `WSDL.Proxy` doesn't coerce datatypes for you. If it did, you could end up calling a completely different function! Good luck debugging that one. It's much easier to be picky about datatypes and fail as quickly as possible if you get them wrong.

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**

```
>>> wsdlFile = 'http://www.xmethods.net/sd/2001/TemperatureService.wsdl'
>>> server y\Client.py", line 453, in __call__      ❶
```

❶

❶

```
<Fault SOAP-ENV:Server:                                        ❷
 Exception from service object: Invalid authorization key: foo:
 <SOAPpy.Types.structType detail at 14164616>:
 {'stackTrace':
  'com.google.soap.search.GoogleSearchFault: Invalid authorization key: foo
    at com.google.soap.search.QueryLimits.lookUpAndLoadFromINSIfNeedBe(
      QueryLimits.java:220)
    at com.google.soap.search.QueryLimits.validateKey(QueryLimits.java:127)
    at com.google.soap.search.GoogleSearchService.doPublicMethodChecks(
      GoogleSearchService.java:825)
    at com.google.soap.search.GoogleSearchService.doGoogleSearch(
      GoogleSearchService.java:121)
    at sun.reflect.GeneratedMethodAccessor13.invoke(Unknown Source)
    at sun.reflect.DelegatingMethodAccessorImpl.invoke(Unknown Source)
    at java.lang.reflect.Method.invoke(Unknown Source)
    at org.apache.soap.server.RPCRouter.invoke(RPCRouter.java:146)
    at org.apache.soap.providers.RPCJavaProvider.invoke(
      RPCJavaProvider.java:129)
    at org.apache.soap.server.http.RPCRouterServlet.doPost(
      RPCRouterServlet.java:288)
    at javax.servlet.http.HttpServlet.service(HttpServlet.java:760)
    at javax.servlet.http.HttpServlet.service(HttpServlet.java:853)
    at com.google.gse.HttpConnection.runServlet(HttpConnection.java:237)
    at com.google.gse.HttpConnection.run(HttpConnection.java:195)
    at com.google.gse.DispatchQueue$WorkerThread.run(DispatchQueue.java:201)
Caused by: com.google.soap.search.UserKeyInvalidException: Key was of wrong size.
    at com.google.soap.search.UserKey.<init>(UserKey.java:59)
    at com.google.soap.search.QueryLimits.lookUpAndLoadFromINSIfNeedBe(
      QueryLimits.java:217)
    ... 14 more
'}>
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
  File "c:\python23\Lib\site-packages\SOAPpy\Client.py", line 453, in __call__
    return self.__r_call(*args, **kw)
  n>", line 1, in ?
  File "c:\python23\Lib\site-packi in ?
```

```
   at com.google.gse.DispatchQueue$WorkerThread.run(DispatchQueue.java:201)
Caused by: com.google.soap.search.UserKeyInvalidException: Key was of wrong size.
   at com.google.soap.search.UserKey.<init>(UserKey.java:59)
   at com.google.soap.search.QueryLimits.lookUpAndLoadFromINSIfNeedBe(
     QueryLimits.java:217)
   ... 14 more
'}>
```

❶   Can you spot the mistake? There's nothing wrong with the calling syntax, or the number of arguments, or the
    datatypes. The problem is application–specific: the first argument is supposed to be my application key, but
    `foo` is not a valid Google key.

❷   The Google server responds with a SOAP Fault and an incredibly long error message, which includes a
    complete Java stack trace. Remember that *all* SOAP errors are signified by SOAP Faults: errors in
    configuration, errors in function arguments, and application–specific errors like this. Buried in there
    somewhere is the crucial piece of information: `Invalid authorization key: foo`.

**Further Reading on Troubleshooting SOAP**

# 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 Python under your belt, you're going to step back and look at the steps that happen *before* the code gets written.

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 mechanics of constructing and validating Roman numerals in Section 7.3, Case Study: Roman Numerals , but now let's step back and consider what it would take to expand that into a two−way utility.

The rules for Roman numerals lead to a number of interesting observations:

1. There is only one correct way to represent a particular number as Roman numerals.
2. 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 way).
3. There is a limited range of numbers that can be expressed as Roman numerals, specifically `1` through `3999`. (The Romans did have several ways of expressing larger numbers, for instance by having a bar over a numeral to represent that its normal value should be multiplied by `1000`, but you're not going to deal with that. For the purposes of this chapter, let's stipulate that Roman numerals go from `1` to `3999`.)
4. There is no way to represent `0` in Roman numerals. (Amazingly, the ancient Romans had no concept of `0` as a number. Numbers were for counting things you had; how can you count what you don't have?)
5. There is no way to represent negative numbers in Roman numerals.
6. There is no way to represent fractions or non−integer numbers in Roman numerals.

Given all of this, what would you expect out of a set of functions to convert to and from Roman numerals?

**`roman.py` requirements**

1. `toRoman` should return the Roman numeral representation for all integers `1` to `3999`.
2. `toRoman` should fail when given an integer outside the range `1` to `3999`.
3. `toRoman` should fail when given a non−integer number.
4. `fromRoman` should take a valid Roman numeral and return the number that it represents.
5. `fromRoman` should fail when given an invalid Roman numeral.
6. 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. So `fromRoman(toRoman(n)) == n` for all `n` in `1..3999`.
7. `toRoman` should always return a Roman numeral using uppercase letters.
8. `fromRoman` should only accept uppercase Roman numerals (*i.e.* it should fail when given lowercase input).

**Further reading**

- This site (http://www.wilkiecollins.demon.co.uk/roman/front.htm) has more on Roman numerals, including a fascinating history (http://www.wilkiecollins.demon.co.uk/roman/intro.htm) of how Romans and other civilizations really used them (short answer: haphazardly and inconsistently).

## 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: you're going to write a test suite that puts these functions through their paces and makes sure that they behave the way you want them to. You read that right: you're going to write code that tests code that you haven't written yet.

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 may become part of later. Python has a framework for unit testing, the appropriately–named `unittest` module.

`unittest` is included with Python 2.1 and later. Python 2.0 users can download it from `pyunit.sourceforge.net` (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 (preferably before writing the code that they test), and to keep them updated as code and requirements change. Unit testing is not a replacement for higher–level functional or system testing, but it is important in all phases of development:

- Before writing code, it forces you to detail your requirements in a useful fashion.
- teet o, ttu m. Td thhe , nobody gomaaoft ooofulr inou to deveriting c o, thhoaveant

```python
class ToRomanBadInput(unittest.TestCase):
    def testTooLarge(self):
        """toRoman should fail with large input"""
        self.assertRaises(roman.OutOfRangeError, roman.toRoman, 4000)

    def testZero(self):
        """toRoman should fail with 0 input"""
        self.assertRaises(roman.OutOfRangeError, roman.toRoman, 0)

    def testNegative(self):
        """toRoman should fail with negative input"""
        self.assertRaises(roman.OutOfRangeError, roman.toRoman, -1)

    def testNonInteger(self):
        """toRoman should fail with non-integer input"""
        self.assertRaises(roman.NotIntegerError, roman.toRoman, 0.5)

class FromRomanBadInput(unittest.TestCase):
    def testTooManyRepeatedNumerals(self):
        """fromRoman should fail with too many repeated numerals"""
        for s in ('MMMM', 'DD', 'CCCC', 'LL', 'XXXX', 'VV', 'IIII'):
            self.assertRaises(roman.InvalidRomanNumeralError, roman.fromRoman, s)

    def testRepeatedPairs(self):
        """fromRoman should fail with repeated pairs of numerals"""
        for s in ('CMCM', 'CDCD', 'XCXC', 'XLXL', 'IXIX', 'IVIV'):
            self.assertRaises(roman.InvalidRomanNumeralError, roman.fromRoman, s)

    def testMalformedAntecedent(self):
        """fromRoman should fail with malformed antecedents"""
        for s in ('IIMXCC', 'VX', 'DCM', 'CMM', 'IXIV',
                  'MCMC', 'XCX', 'IVI', 'LM', 'LD', 'LC'):
            self.assertRaises(roman.InvalidRomanNumeralError, roman.fromRoman, s)

class SanityCheck(unittest.TestCase):
    def testSanity(self):
        """fromRoman(toRoman(n))==n for all n"""
        for integer in range(1, 4000):
            numeral = roman.toRoman(integer)
            result = roman.fromRoman(numeral)
            self.assertEqual(integer, result)

class CaseCheck(unittest.TestCase):
    def testToRomanCase(self):
        """toRoman should always return uppercase"""
        for integer in range(1, 4000):
            numeral = roman.toRoman(integer)
            self.assertEqual(numeral, numeral.upper())

    def testFromRomanCase(self):
        """fromRoman should only accept uppercase input"""
        for integer in range(1, 4000):
            numeral = roman.toRoman(integer)
            roman.fromRoman(numeral.upper())
            self.assertRaises(roman.InvalidRomanNumeralError,
                              roman.fromRoman, numeral.lower())

if __name__ == "__main__":
    unittest.main()
```

**Further reading**

- The PyUnit home page (http://pyunit.sourceforge.net/) has an in–depth discussion of using the `unittest` framework (http://pyunit.sourceforge.net/pyunit.html), including advanced features not covered in this chapter.
- The PyUnit FAQ (http://pyunit.sourceforge.net/pyunit.html) explains why test cases are stored separately (http://pyunit.sourceforge.net/pyunit.html#WHERE) from the code they test.
- *Python Library Reference* (http://www.python.org/doc/current/lib/) summarizes the `unittest` (http://www.python.org/doc/current/lib/module–unittest.html) module.
- ExtremeProgramming.org (http://www.extremeprogramming.org/) discusses why you should write unit tests (http://www.extremeprogramming.org/rules/unittests.html).
- The Portland Pattern Repository (http://www.c2.com/cgi/wiki) has an ongoing discussion of unit tests (http://www.c2.com/cgi/wiki?UnitTests), including a standard definition (http://www.c2.com/cgi/wiki?StandardDefinitionOfUnitTest), why you should code unit tests first (http://www.c2.com/cgi/wiki?CodeUnitTestFirst), and several in–depth case studies (http://www.c2.com/cgi/wiki?UnitTestTrial).

## 13.4. Testing for success

The most fundamental part of unit testing is constructing individual test cases. A test case answers a single question

❶

❷

❸

❹ ❺

❻    Assuming the `toRoman` function was defined correctly, called correctly, completed successfully, and returned a value, the last step is to check whether it returned the *right* value. This is a common question, and the `TestCase` class provides a method, `assertEqual`, to check whether two values are equal. If the result returned from `toRoman` (`result`) does not match the known value you were expecting (`numeral`), `assertEqual` will raise an exception and the test will fail. If the two values are equal, `assertEqual` will do nothing. If every value returned from `toRoman` matches the known value you expect, `assertEqual` never raises an exception, so `testToRomanKnownValues` eventually exits normally, which means `toRoman` has passed this test.

## 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 sort of failure; they must fail in the way you expect.

Remember the other requirements for `toRoman`:

> 2. `toRoman` should fail when given an integer outside the range `1` to `3999`.
> 3. `toRoman` 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 when given bad input.

**Example 13.3. Testing bad input to `toRoman`**

```
class ToRomanBadInput(unittest.TestCase):
    def testTooLarge(self):
        """toRoman should fail with large input"""
        self.assertRaises(roman.OutOfRangeError, roman.toRoman, 4000)  ❶

    def testZero(self):
        """toRoman should fail with 0 input"""
        self.assertRaises(roman.OutOfRangeError, roman.toRoman, 0)     ❷

    def testNegative(self):
        """toRoman should fail with negative input"""
        self.assertRaises(roman.OutOfRangeError, roman.toRoman, -1)

    def testNonInteger(self):
        """toRoman should fail with non-integer input"""
        self.assertRaises(roman.NotIntegerError, roman.toRoman, 0.5)   ❸
```

❶    The `TestCase` class of the `unittest` provides the `assertRaises` method, which takes the following arguments: the exception you're expecting, the function you're testing, and the arguments you're passing that function. (If the function you're testing takes more than one argument, pass them all to `assertRaises`, in order, and it will pass them right along to the function you're testing.) Pay close attention to what you're doing here: instead of calling `toRoman` directly and manually checking that it raises a particular exception (by wrapping it in a `try...except` block), `assertRaises` has encapsulated all of that for us. All you do is give it the exception (`roman.OutOfRangeError`), the function (`toRoman`), and `toRoman`'s arguments (`4000`), and `assertRaises` takes care of calling `toRoman` and checking to make sure that it raises `roman.OutOfRangeError`. (Also note that you're passing the `toRoman` function itself as an argument; you're not calling it, and you're not passing the name of it as a string. Have I mentioned recently how handy it is that everything in Python is an object, including functions and exceptions?)

❷   Along with testing numbers that are too large, you need to test numbers that are too small. Remember, Roman numerals cannot express 0 or negative numbers, so you have a test case for each of those (testZero and testNegative). In testZero, you are testing that toRoman raises a roman.OutOfRangeError exception when called with 0; if it does *not* raise a roman.OutOfRangeError (either because it returns an actual value, or because it raises some other exception), this test is considered failed.

❸   Requirement #3 specifies that toRoman cannot accept a non–integer number, so here you test to make sure that toRoman raises a roman.NotIntegerError exception when called with 0.5. If toRoman does not raise a roman.NotIntegerError, this test is considered failed.

The next two requirements are similar to the first three, except they apply to fromRoman instead of toRoman:

   4. fromRoman should take a valid Roman numeral and return the number that it represents.
   5. fromRoman should fail when given an invalid Roman numeral.

Requirement #4 is handled in the same way as requirement #1, iterating through a sampling of known values and testing each in turn. Requirement #5 is handled in the same way as requirements #2 and #3, by testing a series of bad inputs and making sure fromRoman raises the appropriate exception.

**Example 13.4. Testing bad input to `fromRoman`**

```
class FromRomanBadInput(unittest.TestCase):
    def testTooManyRepeatedNumerals(self):
        """fromRoman should fail with too many repeated numerals"""
        for s in ('MMMM', 'DD', 'CCCC', 'LL', 'XXXX', 'VV', 'IIII'):
            self.assertRaises(roman.InvalidRomanNumeralError, roman.fromRoman, s)  ❶

    def testRepeatedPairs(self):
        """fromRoman should fail with repeated pairs of numerals"""
        for s in ('CMCM', 'CDCD', 'XCXC', 'XLXL', 'IXIX', 'IVIV'):
            self.assertRaises(roman.InvalidRomanNumeralError, roman.fromRoman, s)

    def testMalformedAntecedent(self):
        """fromRoman should fail with malformed antecedents"""
        for s in ('IIMXCC', 'VX', 'DCM', 'CMM', 'IXIV',
                  'MCMC', 'XCX', 'IVI', 'LM', 'LD', 'LC'):
            self.assertRaises(roman.InvalidRomanNumeralError, roman.fromRoman, s)
```

❶   Not much new to say about these; the pattern is exactly the same as the one you used to test bad input to toRoman. I will briefly note that you have another exception: roman.InvalidRomanNumeralError. That makes a total of three custom exceptions that will need to be defined in roman.py (along with roman.OutOfRangeError and roman.NotIntegerError). You'll see how to define these custom exceptions when you actually start writing roman.py, later in this chapter.

# 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 A to B and the other converts B to A. In these cases, it is useful to create a "sanity

6. 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. So `fromRoman(toRoman(n)) == n` for all `n` in `1..3999`.

**Example 13.5. Testing `toRoman` against `fromRoman`**

```
class SanityCheck(unittest.TestCase):
    def testSanity(self):
        """fromRoman(toRoman(n))==n for all n"""
        for integer in range(1, 4000):          ❶ ❷
            numeral = roman.toRoman(integer)
            result = roman.fromRoman(numeral)
            self.assertEqual(integer, result)   ❸
```

❶  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 (`1`) and counting consecutively up to *but not including* the second argument (`4000`). Thus, `1..3999`, which is the valid range for converting to Roman numerals.

❷  I just wanted to mention in passing that `integer` is not a keyword in Python; here it's just a variable name like any other.

❸  The actual testing logic here is straightforward: take a number (`integer`), convert it to a Roman numeral (`numeral`), then convert it back to a number (`result`) and make sure you end up with the same number you started with. If not, `assertEqual` will raise an exception and the test will immediately be considered failed. If all the numbers match, `assertEqual` will always return silently, the entire `testSanity` method will eventually return silently, and the test will be considered passed.

The last two requirements are different from the others because they seem both arbitrary and trivial:

7. `toRoman` should always return a Roman numeral using uppercase letters.
8. `fromRoman` should only accept uppercase Roman numerals (*i.e.* it should fail when given lowercase input).

In fact, they are somewhat arbitrary. You could, for instance, have stipulated that `fromRoman` accept lowercase and mixed case input. But they are not completely arbitrary; if `toRoman` is always returning uppercase output, then `fromRoman` must at least accept uppercase input, or the "sanity check" (requirement #6) would fail. The fact that it *only* accepts uppercase input is arbitrary, but as any systems integrator will tell you, case always matters, so it's worth specifying the behavior up front. And if it's worth specifying, it's worth testing.

**Example 13.6. Testing for case**

```
class CaseCheck(unittest.TestCase):
    def testToRomanCase(self):
        """toRoman should always return uppercase"""
        for integer in range(1, 4000):
            numeral = roman.toRoman(integer)
```

❶ The most interesting thing about this test case is all the things it doesn't test. It doesn't test that the value returned from `toRoman` is right or even consistent; those questions are answered by separate test cases. You have a whole test case just to test for uppercase–ness. You might be tempted to combine this with the sanity check, since both run through the entire range of values and call `toRoman`.[6] But that would violate one of the fundamental rules: each test case should answer only a single question. Imagine that you combined this case check with the sanity check, and then that test case failed. You would need to do further analysis to figure out which part of the test case failed to determine what the problem was. If you need to analyze the results of your unit testing just to figure out what they mean, it's a sure sign that you've mis–designed your test cases.

❷ There's a similar lesson to be learned here: even though "you know" that `toRoman` always returns uppercase, you are explicitly converting its return value to uppercase here to test that `fromRoman` accepts uppercase input. Why? Because the fact that `toRoman` always returns uppercase is an independent requirement. If you changed that requirement so that, for instance, it always returned lowercase, the `testToRomanCase` test case would need to change, but this test case would still work. This was another of the fundamental rules: each test case must be able to work in isolation from any of the others. Every test case is an island.

❸ Note that you're not assigning the return value of `fromRoman` to anything. This is legal syntax in Python; if a function returns a value but nobody's listening, Python just throws away the return value. In this case, that's what you want. This test case doesn't test anything about the return value; it just tests that `fromRoman` accepts the uppercase input without raising an exception.

❹ This is a complicated line, but it's very similar to what you did in the `ToRomanBadInput` and `FromRomanBadInput` tests. You are testing to make sure that calling a particular function (`roman.fromRoman`) with a particular value (`numeral.lower()`, the lowercase version of the current Roman numeral in the loop) raises a particular exception (`roman.InvalidRomanNumeralError`). If it does (each time through the loop), the test passes; if even one time it does something else (like raises a diqaC Ninede8n9

# 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 going to do this in stages, so you can see all the unit tests fail, then watch them pass one by one as you fill in the gaps in `roman.py`.

**Example 14.1. `roman1.py`**

This file is available in `py/roman/stage1/` in the examples directory.

If you have not already done so, you can download this and other examples (http://diveintopython.org/download/diveintopython–examples–5.4.zip) used in this book.

```
"""Convert to and from Roman numerals"""

#Define exceptions
class RomanError(Exception): pass                    ❶
class OutOfRangeError(RomanError): pass              ❷
class NotIntegerError(RomanError): pass
class InvalidRomanNumeralError(RomanError): pass     ❸

def toRoman(n):
    """convert integer to Roman numeral"""
    pass                                             ❹

def fromRoman(s):
    """convert Roman numeral to integer"""
    pass
```

❶    This is how you define your own custom exceptions in Python. Exceptions are classes, and you create your own by subclassing existing exceptions. It is strongly recommended (but not required) that you subclass `Exception`, which is the base class that all built–in exceptions inherit from. Here I am defining `RomanError` (inherited from `Exception`) to act as the base class for all my other custom exceptions to follow. This is a matter of style; I could just as easily have inherited each individual exception from the `Exception` class directly.

❷    The `OutOfRangeError` and `NotIntegerError` exceptions will eventually be used by `toRoman` to flag various forms of invalid input, as specified in `ToRomanBadInput`.

❸    The `InvalidRomanNumeralError` exception will eventually be used by `fromRoman` to flag invalid input, as specified in `FromRomanBadInput`.

❹    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 Python reserved word `pass`.

Now for the big moment (drum roll please): you're finally going to run the unit test against this stubby little module. At this point, every test case should fail. In fact, if any test case passes in stage 1, you should go back to `romantest.py` and re–evaluate why you coded a test so useless that it passes with do–nothing functions.

Run `romantest1.py` 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 any luck, your output should look like this:

**Example 14.2. Output of `romantest1.py` against `roman1.py`**

```
fromRoman should only accept uppercase input ... ERROR
toRoman should always return uppercase ... ERROR
fromRoman should fail with malformed antecedents ... FAIL
fromRoman should fail with repeated pairs of numerals ... FAIL
fromRoman should fail with too many repeated numerals ... FAIL
fromRoman should give known result with known input ... FAIL
toRoman should give known result with known input ... FAIL
fromRoman(toRoman(n))==n for all n ... FAIL
toRoman should fail with non-integer input ... FAIL
toRoman should fail with negative input ... FAIL
toRoman should fail with large input ... FAIL
toRoman should fail with 0 input ... FAIL


======================================================================
ERROR: fromRoman should only accept uppercase input
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage1\romantest1.py", line 154, in testFromRomanCase
    roman1.fromRoman(numeral.upper())
AttributeError: 'None' object has no attribute 'upper'
======================================================================
ERROR: toRoman should always return uppercase
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage1\romantest1.py", line 148, in testToRomanCase
    self.assertEqual(numeral, numeral.upper())
AttributeError: 'None' object has no attribute 'upper'
======================================================================
FAIL: fromRoman should fail with malformed antecedents
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage1\romantest1.py", line 133, in testMalformedAntecedent
    self.assertRaises(roman1.InvalidRomanNumeralError, roman1.fromRoman, s)
  File "c:\python21\lib\unittest.py", line 266, in failUnlessRaises
    raise self.failureException, excName
AssertionError: InvalidRomanNumeralError
======================================================================
FAIL: fromRoman should fail with repeated pairs of numerals
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage1\romantest1.py", line 127, in testRepeatedPairs
    self.assertRaises(roman1.InvalidRomanNumeralError, roman1.fromRoman, s)
  File "c:\python21\lib\unittest.py", line 266, in failUnlessRaises
    raise self.failureException, excName
AssertionError: InvalidRomanNumeralError
======================================================================
FAIL: fromRoman should fail with too many repeated numerals
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage1\romantest1.py", line 122, in testTooManyRepeatedNumerals
    self.assertRaises(roman1.InvalidRomanNumeralError, roman1.fromRoman, s)
  File "c:\python21\lib\unittest.py", line 266, in failUnlessRaises
    raise self.failureException, excName
AssertionError: InvalidRomanNumeralError
======================================================================
FAIL: fromRoman should give known result with known input
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage1\romantest1.py", line 99, in testFromRomanKnownValues
    self.assertEqual(integer, result)
  File "c:\python21\lib\unittest.py", line 273, in failUnlessEqual
    raise self.failureException, (msg or '%s != %s' % (first, second))
AssertionError: 1 != None
```

```
======================================================================
FAIL: toRoman should give known result with known input
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage1\romantest1.py", line 93, in testToRomanKnownValues
    self.assertEqual(numeral, result)
  File "c:\python21\lib\unittest.py", line 273, in failUnlessEqual
```

❶

❷

❸

❹

❶

❷   For each failed test case, `unittest` displays the trace information showing exactly what happened. In this case, the call to `assertRaises` (also called `failUnlessRaises`) raised an `AssertionError` because it was expecting `toRoman` to raise an `OutOfRangeError` and it didn't.

❸   After the detail, `unittest` displays a summary of how many tests were performed and how long it took.

❹   Overall, the unit test failed because at least one test case did not pass. When a test case doesn't pass, `unittest` distinguishes between failures and errors. A failure is a call to an `assertXYZ` method, like `assertEqual` or `assertRaises`, that fails because the asserted condition is not true or the expected exception was not raised. An error is any other sort of exception raised in the code you're testing or the unit test case itself. For instance, the `testFromRomanCase` method ("fromRoman should only accept uppercase input") was an error, because the call to `numeral.upper()` raised an `AttributeError` exception, because `toRoman` was supposed to return a string but didn't. But `testZero` ("toRoman should fail with 0 input") was a failure, because the call to `fromRoman` did not raise the `InvalidRomanNumeral` exception that `assertRaises` was looking for.

## 14.2. `roman.py`, stage 2

Now that you have the framework of the `roman` module laid out, it's time to start writing code and passing test cases.

**Example 14.3. `roman2.py`**

This file is available in `py/roman/stage2/` in the examples directory.

If you have not already done so, you can download this and other examples (http://diveintopython.org/download/diveintopython–examples–5.4.zip) used in this book.

```
"""Convert to and from Roman numerals"""

#Define exceptions
class RomanError(Exception): pass
class OutOfRangeError(RomanError): pass
class NotIntegerError(RomanError): pass
class InvalidRomanNumeralError(RomanError): pass

#Define digit mapping
romanNumeralMap = (('M',  1000),  ❶
                   ('CM', 900),
                   ('D',  500),
                   ('CD', 400),
                   ('C',  100),
                   ('XC', 90),
```
Example,wntopython.org/download/diveintopython–exarMfthis book.

❷

```
def fromRoman(s):
    """convert Roman numeral to integer"""
    pass
```

❶     `romanNumeralMap` is a tuple of tuples which defines three things:

> 1. The character representations of the most basic Roman numerals. Note that this is not just the single–character Roman numerals; you're also defining two–character pairs like `CM` ("one hundred less than one thousand"); this will make the `toRoman` code simpler later.
> 2. The order of the Roman numerals. They are listed in descending value order, from `M` all the way down to `I`.
> 3. The value of each Roman numeral. Each inner tuple is a pair of (*numeral, value*).

❷     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 numerals, you simply iterate through `romanNumeralMap` looking for the largest integer value less than or equal to the input. Once found, you add the Roman numeral representation to the end of the output, subtract the corresponding integer value from the input, lather, rinse, repeat.

### Example 14.4. How `toRoman` works

If you're not clear how `toRoman` works, add a `print` statement to the end of the `while` loop:

```
        while n >= integer:
            result += numeral
            n -= integer
            print 'subtracting', integer, 'from input, adding', numeral, 'to output'
```

```
>>> import roman2
>>> roman2.toRoman(1424)
subtracting 1000 from input, adding M to output
subtracting 400 from input, adding CD to output
subtracting 10 from input, adding X to output
subtracting 10 from input, adding X to output
subtracting 4 from input, adding IV to output
'MCDXXIV'
```

So `toRoman` 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 `romantest2.py` against `roman2.py`

Remember to run `romantest2.py` with the `−v` command–line flag to enable verbose mode.

```
fromRoman should only accept uppercase input ... FAIL
toRoman should always return uppercase ... ok                    ❶
fromRoman should fail with malformed antecedents ... FAIL
fromRoman should fail with repeated pairs of numerals ... FAIL
fromRoman should fail with too many repeated numerals ... FAIL
fromRoman should give known result with known input ... FAIL
toRoman should give known result with known input ... ok         ❷
fromRoman(toRoman(n))==n for all n ... FAIL
toRoman should fail with non-integer input ... FAIL              ❸
toRoman should fail with negative input ... FAIL
toRoman should fail with large input ... FAIL
toRoman should fail with 0 input ... FAIL
```

❶ `toRoman` does, in fact, always return uppercase, because `romanNumeralMap` defines the Roman numeral representations as uppercase. So this test passes already.

❷ Here's the big news: this version of the `toRoman` function passes the known values test. Remember, it's not comprehensive, but it does put the function through its paces with a variety of good inputs, including inputs that produce every single–character Roman numeral, the largest possible input (3999), and the input that produces the longest possible Roman numeral (3888). At this point, you can be reasonably confident that the function works for any good input value you could throw at it.

❸ 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 checks for bad input. Those test cases look for specific exceptions to be raised (via `assertRaises`), and you're never raising them. You'll do that in the next stage.

Here's the rest of the output of the unit test, listing the details of all the failures. You're down to 10.

```
======================================================================
FAIL: fromRoman should only accept uppercase input
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage2\romantest2.py", line 156, in testFromRomanCase
    roman2.fromRoman, numeral.lower())
  File "c:\python21\lib\unittest.py", line 266, in failUnlessRaises
    raise self.failureException, excName
AssertionError: InvalidRomanNumeralError
======================================================================
FAIL: fromRoman should fail with malformed antecedents
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage2\romantest2.py", line 133, in testMalformedAntecedent
    self.assertRaises(roman2.InvalidRomanNumeralError, roman2.fromRoman, s)
  File "c:\python21\lib\unittest.py", line 266, in failUnlessRaises
    raise self.failureException, excName
AssertionError: InvalidRomanNumeralError
======================================================================
FAIL: fromRoman should fail with repeated pairs of numerals
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage2\romantest2.py", line 127, in testRepeatedPairs
    self.assertRaises(roman2.InvalidRomanNumeralError, roman2.fromRoman, s)
  File "c:\python21\lib\unittest.py", line 266, in failUnlessRaises
    raise self.failureException, excName
AssertionError: InvalidRomanNumeralError
======================================================================
FAIL: fromRoman should fail with too many repeated numerals
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage2\romantest2.py", line 122, in testTooManyRepeatedNumerals
    self.assertRaises(roman2.InvalidRomanNumeralError, roman2.fromRoman, s)
  File "c:\python21\lib\unittest.py", line 266, in failUnlessRaises
    raise self.failureException, excName
AssertionError: InvalidRomanNumeralError
======================================================================
FAIL: fromRoman should give known result with known input
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage2\romantest2.py", line 99, in testFromRomanKnownValues
    self.assertEqual(integer, result)
  File "c:\python21\lib\unittest.py", line 273, in failUnlessEqual
    raise self.failureException, (msg or '%s != %s' % (first, second))
AssertionError: 1 != None
======================================================================
FAIL: fromRoman(toRoman(n))==n for all n
----------------------------------------------------------------------
```

```
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage2\romantest2.py", line 141, in testSanity
    self.assertEqual(integer, result)
  File "c:\python21\lib\unittest.py", line 273, in failUnlessEqual
    raise self.failureException, (msg or '%s != %s' % (first, second))
AssertionError: 1 != None
======================================================================
FAIL: toRoman should fail with non-integer input
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage2\romantest2.py", line 116, in testNonInteger
    self.assertRaises(roman2.NotIntegerError, roman2.toRoman, 0.5)
  File "c:\python21\lib\unittest.py", line 266, in failUnlessRaises
    raise self.failureException, excName
AssertionError: NotIntegerError
======================================================================
FAIL: toRoman should fail with negative input
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage2\romantest2.py", line 112, in testNegative
    self.assertRaises(roman2.OutOfRangeError, roman2.toRoman, -1)
  File "c:\python21\lib\unittest.py", line 266, in failUnlessRaises
    raise self.failureException, excName
AssertionError: OutOfRangeError
======================================================================
FAIL: toRoman should fail with large input
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage2\romantest2.py", line 104, in testTooLarge
    self.assertRaises(roman2.OutOfRangeError, roman2.toRoman, 4000)
  File "c:\python21\lib\unittest.py", line 266, in failUnlessRaises
    raise self.failureException, excName
AssertionError: OutOfRangeError
======================================================================
FAIL: toRoman should fail with 0 input
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage2\romantest2.py", line 108, in testZero
    self.assertRaises(roman2.OutOfRangeError, roman2.toRoman, 0)
  File "c:\python21\lib\unittest.py", line 266, in failUnlessRaises
    raise self.failureException, excName
AssertionError: OutOfRangeError
----------------------------------------------------------------------
Ran 12 tests in 0.320s

FAILED (failures=10)
```

## 14.3. `roman.py`, stage 3

Now that `toRoman` behaves correctly with good input (integers from `1` to `3999`), it's time to make it behave correctly with bad input (everything else).

**Example 14.6. `roman3.py`**

This file is available in `py/roman/stage3/` in the examples directory.

If you have not already done so, you can download this and other examples (http://diveintopython.org/download/diveintopython−examples−5.4.zip) used in this book.

```
"""Convert to and from Roman numerals"""

#Define exceptions
class RomanError(Exception): pass
class OutOfRangeError(RomanError): pass
class NotIntegerError(RomanError): pass
class InvalidRomanNumeralError(RomanError): pass

#Define digit mapping
romanNumeralMap = (('M',  1000),
                   ('CM', 900),
                   ('D',  500),
                   ('CD', 400),
                   ('C',  100),
                   ('XC', 90),
                   ('L',  50),
                   ('XL', 40),
                   ('X',  10),
                   ('IX', 9),
                   ('V',  5),
                   ('IV', 4),
                   ('I',  1))

def toRoman(n):
    """convert integer to Roman numeral"""
    if not (0 < n < 4000):                                          ❶
        raise OutOfRangeError, "number out of range (must be 1..3999)" ❷
    if int(n) <> n:                                                 ❸
        raise NotIntegerError, "non-integers can not be converted"

    result = ""                                                     ❹
    for numeral, integer in romanNumeralMap:
        while n >= integer:
            result += numeral
            n -= integer
    return result

def fromRoman(s):
    """convert Roman numeral to integer"""
    pass
```

❶  This is a nice Pythonic shortcut: multiple comparisons at once. This is equivalent to `if not ((0 < n)` `and (n < 4000))`, but it's much easier to read. This is the range check, and it should catch inputs that are too large, negative, or zero.

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 that you've defined. The second parameter, the error message, is optional; if given, it is displayed in the traceback that is printed if the exception is never handled.

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 `toRoman` handle bad input**

```
>>> import roman3
>>> roman3.toRoman(4000)
Traceback (most recent call last):
  File "<interactive input>", line 1, in ?
  File "roman3.py", line 27, in toRoman
    raise OutOfRangeError, "number out of range (must be 1..3999)"
OutOfRangeError: number out of range (must be 1..3999)
>>> roman3.toRoman(1.5)
```

```
Traceback (most recent call last):
  File "<interactive input>", line 1, in ?
  File "roman3.py", line 29, in toRoman
    raise NotIntegerError, "non-integers can not be converted"
NotIntegerError: non-integers can not be converted
```

**Example 14.8. Output of `romantest3.py` against `roman3.py`**

```
fromRoman should only accept uppercase input ... FAIL
toRoman should always return uppercase ... ok
fromRoman should fail with malformed antecedents ... FAIL
fromRoman should fail with repeated pairs of numerals ... FAIL
fromRoman should fail with too many repeated numerals ... FAIL
fromRoman should give known result with known input ... FAIL
toRoman should give known result with known input ... ok ❶
fromRoman(toRoman(n))==n for all n ... FAIL
toRoman should fail with non-integer input ... ok        ❷
toRoman should fail with negative input ... ok           ❸
toRoman should fail with large input ... ok
toRoman should fail with 0 input ... ok
```

❶   `toRoman` 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.

❷   More exciting is the fact that all of the bad input tests now pass. This test, `testNonInteger`, passes because of the `int(n) <> n` check. When a non–integer is passed to `toRoman`, the `int(n) <> n` check notices it and raises the `NotIntegerError` exception, which is what `testNonInteger` is looking for.

❸   This test, `testNegative`, passes because of the `not (0 < n < 4000)` check, which raises an `OutOfRangeError` exception, which is what `testNegative` is looking for.

```
======================================================================
FAIL: fromRoman should only accept uppercase input
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage3\romantest3.py", line 156, in testFromRomanCase
    roman3.fromRoman, numeral.lower())
  File "c:\python21\lib\unittest.py", line 266, in failUnlessRaises
```

```
File "C:\docbook\dip\py\roman\stage3\romantest3.py", line 122, in testTooManyRepeatedNumerals
    self.assertRaises(roman3.InvalidRomanNumeralError, roman3.fromRoman, s)
```

❶

❶

☞

```
#Define digit mapping
romanNumeralMap = (('M',  1000),
                   ('CM', 900),
                   ('D',  500),
                   ('CD', 400),
                   ('C',  100),
                   ('XC', 90),
                   ('L',  50),
                   ('XL', 40),
                   ('X',  10),
                   ('IX', 9),
                   ('V',  5),
                   ('IV', 4),
                   ('I',  1))

# toRoman function omitted for clarity (it hasn't changed)

def fromRoman(s):
    """convert Roman numeral to integer"""
    result = 0
    index = 0
    for numeral, integer in romanNumeralMap:
        while s[index:index+len(numeral)] == numeral:    ❶
            result += integer
            index += len(numeral)
    return result
```

❶    The pattern here is the same as toRoman. You iterate through your Roman numeral data structure (a tuple of tuples), and instead of matching the highest integer values as often as possible, you match the "highest" Roman numeral character strings as often as possible.

**Example 14.10. How `fromRoman` works**

If you're not clear how fromRoman works, add a print statement to the end of the while loop:

```
        while s[index:index+len(numeral)] == numeral:
            result += integer
            index += len(numeral)
            print 'found', numeral, 'of length', len(numeral), ', adding', integer
```

```
>>> import roman4
>>> roman4.fromRoman('MCMLXXII')
found M , of length 1, adding 1000
found CM , of length 2, adding 900
found L , of length 1, adding 50
found X , of length 1, adding 10
found X , of length 1, adding 10
found I , of length 1, adding 1
found I , of length 1, adding 1
1972
```

**Example 14.11. Output of `romantest4.py` against `roman4.py`**

```
fromRoman should only accept uppercase input ... FAIL
toRoman should always return uppercase ... ok
fromRoman should fail with malformed antecedents ... FAIL
fromRoman should fail with repeated pairs of numerals ... FAIL
fromRoman should fail with too many repeated numerals ... FAIL
fromRoman should give known result with known input ... ok  ❶
```

```
toRoman should give known result with known input ... ok
fromRoman(toRoman(n))==n for all n ... ok                    ❷
toRoman should fail with non-integer input ... ok
toRoman should fail with negative input ... ok
toRoman should fail with large input ... ok
toRoman should fail with 0 input ... ok
```

❶    Two pieces of exciting news here. The first is that `fromRoman` 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 `toRoman` and `fromRoman` work properly for all possible good values. (This is not guaranteed; it is theoretically possible that `toRoman` has a bug that produces the wrong Roman numeral for some particular set of inputs, *and* that `fromRoman` has a reciprocal bug that produces the same wrong integer values for exactly that set of Roman numerals that `toRoman` generated incorrectly. Depending on your application and your requirements, this possibility may bother you; if so, write more comprehensive test cases until it doesn't bother you.)

```
======================================================================
FAIL: fromRoman should only accept uppercase input
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage4\romantest4.py", line 156, in testFromRomanCase
    roman4.fromRoman, numeral.lower())
  File "c:\python21\lib\unittest.py", line 266, in failUnlessRaises
    raise self.failureException, excName
AssertionError: InvalidRomanNumeralError
======================================================================
FAIL: fromRoman should fail with malformed antecedents
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage4\romantest4.py", line 133, in testMalformedAntecedent
    self.assertRaises(roman4.InvalidRomanNumeralError, roman4.fromRoman, s)
  File "c:\python21\lib\unittest.py", line 266, in failUnlessRaises
    raise self.failureException, excName
AssertionError: InvalidRomanNumeralError
======================================================================
FAIL: fromRoman should fail with repeated pairs of numerals
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage4\romantest4.py", line 127, in testRepeatedPairs
    self.assertRaises(roman4.InvalidRomanNumeralError, roman4.fromRoman, s)
  File "c:\python21\lib\unittest.py", line 266, in failUnlessRaises
    raise self.failureException, excName
AssertionError: InvalidRomanNumeralError
======================================================================
FAIL: fromRoman should fail with too many repeated numerals
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage4\romantest4.py", line 122, in testTooManyRepeatedNumerals
    self.assertRaises(roman4.InvalidRomanNumeralError, roman4.fromRoman, s)
  File "c:\python21\lib\unittest.py", line 266, in failUnlessRaises
    raise self.failureException, excName
AssertionError: InvalidRomanNumeralError
----------------------------------------------------------------------
Ran 12 tests in 1.222s

FAILED (failures=4)
```

# 14.5. `roman.py`, stage 5

Now that `fromRoman` 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 finding a way to look at a string and determine if it's a valid Roman numeral. This is inherently more difficult than validating numeric input in `toRoman`, but you have a powerful tool at your disposal: regular expressions.

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 letters `M`, `D`, `C`, `L`, `X`, `V`, and `I`. Let's review the rules:

1. Characters are additive. `I` is 1, `II` is 2, and `III` is 3. `VI` is 6 (literally, "5 and 1"), `VII` is 7, and `VIII` is 8.
2. The tens characters (`I`, `X`, `C`, and `M`) can be repeated up to three times. At 4, you need to subtract from the next highest fives character. You can't represent 4 as `IIII`; instead, it is represented as `IV` ("1 less than 5"). 40 is written as `XL` ("10 less than 50"), 41 as `XLI`, 42 as `XLII`, 43 as `XLIII`, and then 44 as `XLIV` ("10 less than 50, then 1 less than 5").
3. Similarly, at 9, you need to subtract from the next highest tens character: 8 is `VIII`, but 9 is `IX` ("1 less than 10"), not `VIIII` (since the `I` character can not be repeated four times). 90 is `XC`, 900 is `CM`.
4. The fives characters can not be repeated. 10 is always represented as `X`, never as `VV`. 100 is always `C`, never `LL`.
5. Roman numerals are always written highest to lowest, and read left to right, so order of characters matters very much. `DC` is 600; `CD` is a completely different number (400, "100 less than 500"). `CI` is 101; `IC` is not even a valid Roman numeral (because you can't subtract 1 directly from 100; you would need to write it as `XCIX`, "10 less than 100, then 1 less than 10").

**Example 14.12. `roman5.py`**

This file is available in `py/roman/stage5/` in the examples directory.

If you have not already done so, you can download this and other examples (http://diveintopython.org/download/diveintopython–examples–5.4.zip) used in this book.

```
"""Convert to and from Roman numerals"""
import re

#Define exceptions
class RomanError(Exception): pass
class OutOfRangeError(RomanError): pass
class NotIntegerError(RomanError): pass
class InvalidRomanNumeralError(RomanError): pass

#Define digit mapping
romanNumeralMap = (('M',  1000),
                   ('CM', 900),
                   ('D',  500),
                   ('CD', 400),
                   ('C',  100),
                   ('XC', 90),
                   ('L',  50),
                   ('XL', 40),
                   ('X',  10),
                   ('IX', 9),
                   ('V',  5),
```

❶

❷

❶

❷

❶

❷
❸

OK                                                                    ❹

❶

❷

❸

❹

# 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 yet.

**Example 15.1. The bug**

```
>>> import roman5
>>> roman5.fromRoman("") ❶
0
```

❶    Remember in the previous section when you kept seeing that an empty string would match the regular expression you were using to check for valid Roman numerals? Well, it turns out that this is still true for the final version of the regular expression. And that's a bug; you want an empty string to raise an `InvalidRomanNumeralError` exception just like any other sequence of characters that don't represent a valid Roman numeral.

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 (`romantest61.py`)**

```
class FromRomanBadInput(unittest.TestCase):

    # previous test cases omitted for clarity (they haven't changed)

    def testBlank(self):
        """fromRoman should fail with blank string"""
        self.assertRaises(roman.InvalidRomanNumeralError, roman.fromRoman, "") ❶
```

❶    Pretty simple stuff here. Call `fromRoman` with an empty string and make sure it raises an `InvalidRomanNumeralError` exception. The hard part was finding the bug; now that you know about it, testing for it is the easy part.

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 `romantest61.py` against `roman61.py`**

```
fromRoman should only accept uppercase input ... ok
toRoman should always return uppercase ... ok
fromRoman should fail with blank string ... FAIL
fromRoman should fail with malformed antecedents ... ok
fromRoman should fail with repeated pairs of numerals ... ok
fromRoman should fail with too many repeated numerals ... ok
fromRoman should give known result with known input ... ok
toRoman should give known result with known input ... ok
fromRoman(toRoman(n))==n for all n ... ok
toRoman should fail with non-integer input ... ok
toRoman should fail with negative input ... ok
toRoman should fail with large input ... ok
toRoman should fail with 0 input ... ok

======================================================================
```

```
FAIL: fromRoman should fail with blank string
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage6\romantest61.py", line 137, in testBlank
    self.assertRaises(roman61.InvalidRomanNumeralError, roman61.fromRoman, "")
  File "c:\python21\lib\unittest.py", line 266, in failUnlessRaises
    raise self.failureException, excName
AssertionError: InvalidRomanNumeralError
----------------------------------------------------------------------
Ran 13 tests in 2.864s

FAILED (failures=1)
```

*Now* you can fix the bug.

### Example 15.4. Fixing the bug (`roman62.py`)

This file is available in `py/roman/stage6/` in the examples directory.

```
def fromRoman(s):
    """convert Roman numeral to integer"""
    if not s: ❶
        raise InvalidRomanNumeralError, 'Input can not be blank'
    if not re.search(romanNumeralPattern, s):
        raise InvalidRomanNumeralError, 'Invalid Roman numeral: %s' % s

    result = 0
    index = 0
    for numeral, integer in romanNumeralMap:
        while s[index:index+len(numeral)] == numeral:
            result += integer
            index += len(numeral)
    return result
```

❶    Only two lines of code are required: an explicit check for an empty string, and a `raise` statement.

### Example 15.5. Output of `romantest62.py` against `roman62.py`

```
fromRoman should only accept uppercase input ... ok
toRoman should always return uppercase ... ok
fromRoman should fail with blank string ... ok ❶
fromRoman should fail with malformed antecedents ... ok
fromRoman should fail with repeated pairs of numerals ... ok
fromRoman should fail with too many repeated numerals ... ok
fromRoman should give known result with known input ... ok
toRoman should give known result with known input ... ok
fromRoman(toRoman(n))==n for all n ... ok
toRoman should fail with non-integer input ... ok
toRoman should fail with negative input ... ok
toRoman should fail with large input ... ok
toRoman should fail with 0 input ... ok


----------------------------------------------------------------------
Ran 13 tests in 2.834s

OK ❷
```

❶    The blank string test case now passes, so the bug is fixed.

❷    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. In a testing–centric environment, it may *seem* like it takes longer to fix a bug, since you need to articulate in code exactly what the bug is (to write the test case), then fix the bug itself. Then if the test case doesn't pass right away, you need to figure out whether the fix was wrong, or whether the test case itself has a bug in it. However, in the long run, this back–and–forth between test code and code tested pays for itself, because it makes it more likely that bugs are fixed correctly the first time. Also, since you can easily re–run *all* the test cases along with your new one, you are much less likely to break old code when fixing new code. Today's unit test is tomorrow's regression test.

# 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 involving scissors and hot wax, requirements will change. Most customers don't know what they want until they see it, and even if they do, they aren't that good at articulating what they want precisely enough to be useful. And even if they do, they'll want more in the next release anyway. So be prepared to update your test cases as requirements change.

Suppose, for instance, that you wanted to expand the range of the Roman numeral conversion functions. Remember the rule that said that no character could be repeated more than three times? Well, the Romans were willing to make an exception to that rule by having 4 M characters in a row to represent 4000. If you make this change, you'll be able to expand the range of convertible numbers from 1..3999 to 1..4999. But first, you need to make some changes to the test cases.

**Example 15.6. Modifying test cases for new requirements (`romantest71.py`)**

This file is available in py/roman/stage7/ in the examples directory.

If you have not already done so, you can download this and other examples
(http://diveintopython.org/download/diveintopython–examples–5.4.zip) used in this book.

```
import roman71
import unittest

class KnownValues(unittest.TestCase):
    knownValues = ( (1, 'I'),
                    (2, 'II'),
                    (3, 'III'),
                    (4, 'IV'),
                    (5, 'V'),
                    (6, 'VI'),
                    (7, 'VII'),
                    (8, 'VIII'),
                    (9, 'IX'),
                    (10, 'X'),
                    (50, 'L'),
                    (100, 'C'),
                    (500, 'D'),
                    (1000, 'M'),
                    (31, 'XXXI'),
                    (148, 'CXLVIII'),
                    (294, 'CCXCIV'),
                    (312, 'CCCXII'),
                    (421, 'CDXXI'),
                    (528, 'DXXVIII'),
```

```
                    (621, 'DCXXI'),
                    (782, 'DCCLXXXII'),
                    (870, 'DCCCLXX'),
                    (941, 'CMXLI'),
                    (1043, 'MXLIII'),
                    (1110, 'MCX'),
                    (1226, 'MCCXXVI'),
                    (1301, 'MCCCI'),
                    (1485, 'MCDLXXXV'),
                    (1509, 'MDIX'),
                    (1607, 'MDCVII'),
                    (1754, 'MDCCLIV'),
                    (1832, 'MDCCCXXXII'),
                    (1993, 'MCMXCIII'),
                    (2074, 'MMLXXIV'),
                    (2152, 'MMCLII'),
                    (2212, 'MMCCXII'),
                    (2343, 'MMCCCXLIII'),
                    (2499, 'MMCDXCIX'),
                    (2574, 'MMDLXXIV'),
                    (2646, 'MMDCXLVI'),
                    (2723, 'MMDCCXXIII'),
                    (2892, 'MMDCCCXCII'),
                    (2975, 'MMCMLXXV'),
                    (3051, 'MMMLI'),
                    (3185, 'MMMCLXXXV'),
                    (3250, 'MMMCCL'),
                    (3313, 'MMMCCCXIII'),
                    (3408, 'MMMCDVIII'),
                    (3501, 'MMMDI'),
                    (3610, 'MMMDCX'),
                    (3743, 'MMMDCCXLIII'),
                    (3844, 'MMMDCCCXLIV'),
                    (3888, 'MMMDCCCLXXXVIII'),
                    (3940, 'MMMCMXL'),
                    (3999, 'MMMCMXCIX'),
                    (4000, 'MMMM'),                                    ❶
                    (4500, 'MMMMD'),
                    (4888, 'MMMMDCCCLXXXVIII'),
                    (4999, 'MMMMCMXCIX'))

    def testToRomanKnownValues(self):
        """toRoman should give known result with known input"""
        for integer, numeral in self.knownValues:
            result = roman71.toRoman(integer)
            self.assertEqual(numeral, result)

    def testFromRomanKnownValues(self):
        """fromRoman should give known result with known input"""
        for integer, numeral in self.knownValues:
            result = roman71.fromRoman(numeral)
            self.assertEqual(integer, result)

class ToRomanBadInput(unittest.TestCase):
    def testTooLarge(self):
        """toRoman should fail with large input"""
        self.assertRaises(roman71.OutOfRangeError, roman71.toRoman, 5000)  ❷

    def testZero(self):
        """toRoman should fail with 0 input"""
        self.assertRaises(roman71.OutOfRangeError, roman71.toRoman, 0)

    def testNegative(self):
```

```
        """toRoman should fail with negative input"""
        self.assertRaises(roman71.OutOfRangeError, roman71.toRoman, -1)

    def testNonInteger(self):
        """toRoman should fail with non-integer input"""
        self.assertRaises(roman71.NotIntegerError, roman71.toRoman, 0.5)

class FromRomanBadInput(unittest.TestCase):
    def testTooManyRepeatedNumerals(self):
        """fromRoman should fail with too many repeated numerals"""
        for s in ('MMMM', 'DD', 'CCCC', 'LL', 'XXXX', 'VV', 'IIII'):     ❸
            self.assertRaises(roman71.InvalidRomanNumeralError, roman71.fromRoman, s)

    def testRepeatedPairs(self):
        """fromRoman should fail with repeated pairs of numerals"""
        for s in ('CMCM', 'CDCD', 'XCXC', 'XLXL', 'IXIX', 'IVIV'):
            self.assertRaises(roman71.InvalidRomanNumeralError, roman71.fromRoman, s)

    def testMalformedAntecedent(self):
        """fromRoman should fail with malformed antecedents"""
        for s in ('IIMXCC', 'VX', 'DCM', 'CMM', 'IXIV',
                  'MCMC', 'XCX', 'IVI', 'LM', 'LD', 'LC'):
            self.assertRaises(roman71.InvalidRomanNumeralError, roman71.fromRoman, s)

    def testBlank(self):
        """fromRoman should fail with blank string"""
        self.assertRaises(roman71.InvalidRomanNumeralError, roman71.fromRoman, "")

class SanityCheck(unittest.TestCase):
    def testSanity(self):
        """fromRoman(toRoman(n))==n for all n"""
        for integer in range(1, 5000):                                   ❹
            numeral = roman71.toRoman(integer)
            result = roman71.fromRoman(numeral)
            self.assertEqual(integer, result)

class CaseCheck(unittest.TestCase):
    def testToRomanCase(self):
        """toRoman should always return uppercase"""
        for integer in range(1, 5000):
            numeral = roman71.toRoman(integer)
            self.assertEqual(numeral, numeral.upper())

    def testFromRomanCase(self):
        """fromRoman should only accept uppercase input"""
        for integer in range(1, 5000):
            numeral = roman71.toRoman(integer)
            roman71.fromRoman(numeral.upper())
            self.assertRaises(roman71.InvalidRomanNumeralError,
                              roman71.fromRoman, numeral.lower())

if __name__ == "__main__":
    unittest.main()
```

❶ 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 4000 range. Here I've included 4000 (the shortest), 4500 (the second shortest), 4888 (the longest), and 4999 (the largest).

❷ The definition of "large input" has changed. This test used to call toRoman with 4000 and expect an error; now that 4000-4999 are good values, you need to bump this up to 5000.

❸ The definition of "too many repeated numerals" has also changed. This test used to call fromRoman with 'MMMM' and expect an error; now that MMMM is considered a valid Roman numeral, you need to

bump this up to `'MMMMM'`.

❹ The sanity check and case checks loop through every number in the range, from `1` to `3999`. Since the range has now expanded, these `for` loops need to be updated as well to go up to `4999`.

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 `romantest71.py` against `roman71.py`**

❶

❷
❸
❹

❶

❷

❸

❹

```
ERROR: toRoman should give known result with known input
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage7\romantest71.py", line 96, in testToRomanKnownValues
    result = roman71.toRoman(integer)
  File "roman71.py", line 28, in toRoman
    raise OutOfRangeError, "number out of range (must be 1..3999)"
OutOfRangeError: number out of range (must be 1..3999)
======================================================================
ERROR: fromRoman(toRoman(n))==n for all n
----------------------------------------------------------------------
Traceback (most recent call last):
  File "C:\docbook\dip\py\roman\stage7\romantest71.py", line 147, in testSanity
    numeral = roman71.toRoman(integer)
  File "roman71.py", line 28, in toRoman
    raise OutOfRangeError, "number out of range (must be 1..3999)"
OutOfRangeError: number out of range (must be 1..3999)
----------------------------------------------------------------------
Ran 13 tests in 2.213s

FAILED (errors=5)
```

Now that you have test cases that fail due to the new requirements, you can think about fixing the code to bring it in line with the test cases. (One thing that takes some getting used to when you first start coding unit tests is that the code being tested is never "ahead" of the test cases. While it's behind, you still have some work to do, and as soon as it catches up to the test cases, you stop coding.)

**Example 15.8. Coding the new requirements (`roman72.py`)**

This file is available in `py/roman/stage7/` in the examples directory.

```
"""Convert to and from Roman numerals"""
import re

#Define exceptions
class RomanError(Exception): pass
class OutOfRangeError(RomanError): pass
class NotIntegerError(RomanError): pass
class InvalidRomanNumeralError(RomanError): pass

#Define digit mapping
romanNumeralMap = (('M',  1000),
                   ('CM', 900),
                   ('D',  500),
                   ('CD', 400),
                   ('C',  100),
                   ('XC', 90),
                   ('L',  50),
                   ('XL', 40),
                   ('X',  10),
                   ('IX', 9),
                   ('V',  5),
                   ('IV', 4),
                   ('I',  1))

def toRoman(n):
    """convert integer to Roman numeral"""
    if not (0 < n < 5000):                                                    ❶
        raise OutOfRangeError, "number out of range (must be 1..4999)"
    if int(n) <> n:
```

❷

❶

❷

```
Ran 13 tests in 3.685s

OK ❶
```

❶     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 someone else blames you for breaking their code and you can actually *prove* that you didn't. The best thing about unit testing is that it gives you the freedom to refactor mercilessly.

Refactoring is the process of taking working code and making it work better. Usually, "better" means "faster", although it can also mean "using less memory", or "using less disk space", or simply "more elegantly". Whatever it means to you, to your project, in your environment, refactoring is important to the long–term health of any program.

Here, "better" means "faster". Specifically, the `fromRoman` function is slower than it needs to be, because of that big nasty regular expression that you use to validate Roman numerals. It's probably not worth trying to do away with the regular expression altogether (it would be difficult, and it might not end up any faster), but you can speed up the function by precompiling the regular expression.

**Example 15.10. Compiling regular expressions**

```
>>> import re
>>> pattern = '^M?M?M?$'
>>> re.search(pattern, 'M')                    ❶
<SRE_Match object at 01090490>
>>> compiledPattern = re.compile(pattern) ❷
>>> compiledPattern
<SRE_Pattern object at 00F06E28>
>>> dir(compiledPattern)                       ❸
['findall', 'match', 'scanner', 'search', 'split', 'sub', 'subn']
>>> compiledPattern.search('M')                ❹
<SRE_Match object at 01104928>
```

❶     This is the syntax you've seen before: `re.search` takes a regular expression as a string (`pattern`) and a
    string to match against it (`'M'`). If the pattern matches, the function returns a match object which can be
    queried to find out exactly what matched and how.

❷     This is the new syntax: `re.compile` takes a regular expression as a string and returns a pattern object. Note
    there is no string to match here. Compiling a regular expression has nothing to do with matching it against any
    specific strings (like `'M'`); it only involves the regular expression itself.

❸     The compiled pattern object returned from `re.compile` has several useful–looking functions, including
    several (like `search` and `sub`) that are available directly in the `re` module.

❹     Calling the compiled pattern object's `search` function with the string `'M'` accomplishes the same thing as
    calling `re.search` with both the regular expression and the string `'M'`. Only much, much faster. (In fact, the
    `re.search` function simply compiles the regular expression and calls the resulting pattern object's `search`
    method for you.)

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 directly.

**Example 15.11. Compiled regular expressions in `roman81.py`**

This file is available in `py/roman/stage8/` in the examples directory.

If you have not already done so, you can download this and other examples
(http://diveintopython.org/download/diveintopython–examples–5.4.zip) used in this book.

```
# toRoman and rest of module omitted for clarity

romanNumeralPattern = \
    re.compile('^M?M?M?M?(CM|CD|D?C?C?C?)(XC|XL|L?X?X?X?)(IX|IV|V?I?I?I?)$')  ❶

def fromRoman(s):
    """convert Roman numeral to integer"""
    if not s:
        raise InvalidRomanNumeralError, 'Input can not be blank'
    if not romanNumeralPattern.search(s):                                      ❷
        raise InvalidRomanNumeralError, 'Invalid Roman numeral: %s' % s

    result = 0
    index = 0
    for numeral, integer in romanNumeralMap:
        while s[index:index+len(numeral)] == numeral:
            result += integer
            index += len(numeral)
    return result
```

❶   This looks very similar, but in fact a lot has changed. `romanNumeralPattern` is no longer a string; it is a pattern object which was returned from `re.compile`.

❷   That means that you can call methods on `romanNumeralPattern` directly. This will be much, much faster than calling `re.search` every time. The regular expression is compiled once and stored in `romanNumeralPattern` when the module is first imported; then, every time you call `fromRoman`, you can immediately match the input string against the regular expression, without any intermediate steps occurring under the covers.

So 27Mwmuch faster tf (3a)] a0oing

❶

❷

❸

❶

❷

❸

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 there is frequently more than one way to write the same expression. After some discussion about this module on comp.lang.python (http://groups.google.com/groups?group=comp.lang.python), someone suggested that I try using the {m,n} syntax for the optional repeated characters.

**Example 15.13. `roman82.py`**

This file is available in `py/roman/stage8/` in the examples directory.

If you have not already done so, you can download this and other examples (http://diveintopython.org/download/diveintopython–examples–5.4.zip) used in this book.

```
# rest of program omitted for clarity

#old version
#romanNumeralPattern = \
#   re.compile('^M?M?M?M?(CM|CD|D?C?C?C?)(XC|XL|L?X?X?X?)(IX|IV|V?I?I?I?)$')

#new version
romanNumeralPattern = \
    re.compile('^M{0,4}(CM|CD|D?C{0,3})(XC|XL|L?X{0,3})(IX|IV|V?I{0,3})$')  ❶
```

❶   You have replaced M?M?M?M? with M{0,4}. Both mean the same thing: "match 0 to 4 M characters". Similarly, C?C?C? became C{0,3} ("match 0 to 3 C characters") and so forth for X and I.

This form of the regular expression is a little shorter (though not any more readable). The big question is, is it any faster?

**Example 15.14. Output of `romantest82.py` against `roman82.py`**

```
.............
----------------------------------------------------------------------
Ran 13 tests in 3.315s  ❶

OK                      ❷
```

❶   Overall, the unit tests run 2% faster with this form of regular expression. That doesn't sound exciting, but remember that the `search` function is a small part of the overall unit test; most of the time is spent doing other things. (Separately, I time–tested just the regular expressions, and found that the `search` function is 11% faster with this syntax.) By precompiling the regular expression and rewriting part of it to use this new syntax, you've improved the regular expression performance by over 60%, and improved the overall performance of the entire unit test by over 10%.

❷   More important than any performance boost is the fact that the module still works perfectly. This is the freedom I was talking about earlier: the freedom to tweak, change, or rewrite any piece of it and verify that you haven't messed anything up in the process. This is not a license to endlessly tweak your code just for the sake of tweaking it; you had a very specific objective ("make `fromRoman` faster"), and you were able to accomplish that objective without any lingering doubts about whether you introduced new bugs in the process.

One other tweak I would like to make, and then I promise I'll stop refactoring and put this module to bed. As you've seen repeatedly, regular expressions can get pretty hairy and unreadable pretty quickly. I wouldn't like to come back to this module in six months and try to maintain it. Sure, the test cases pass, so I know that it works, but if I can't figure out *how* it works, it's still going to be difficult to add new features, fix new bugs, or otherwise maintain it. As you saw

❶

❶

❶
❷

❶

❷

you build the lookup table for converting integers to Roman numerals, you can build the reverse lookup table to convert Roman numerals to integers.

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 he could prove that his code worked just as well as the original.

**Example 15.17. `roman9.py`**

This file is available in `py/roman/stage9/` in the examples directory.

If you have not already done so, you can download this and other examples (http://diveintopython.org/download/diveintopython−examples−5.4.zip) used in this book.

```
#Define exceptions
class RomanError(Exception): pass
class OutOfRangeError(RomanError): pass
class NotIntegerError(RomanError): pass
class InvalidRomanNumeralError(RomanError): pass

#Roman numerals must be less than 5000
MAX_ROMAN_NUMERAL = 4999

#Define digit mapping
romanNumeralMap = (('M',  1000),
                   ('CM', 900),
                   ('D',  500),
                   ('CD', 400),
                   ('C',  100),
                   ('XC', 90),
                   ('L',  50),
                   ('XL', 40),
                   ('X',  10),
                   ('IX', 9),
                   ('V',  5),
                   ('IV', 4),
                   ('I',  1))

#Create tables for fast conversion of roman numerals.
#See fillLookupTables() below.
toRomanTable = [ None ]  # Skip an index since Roman numerals have no zero
fromRomanTable = {}

def toRoman(n):
    """convert integer to Roman numeral"""
    if not (0 < n <= MAX_ROMAN_NUMERAL):
        raise OutOfRangeError, "number out of range (must be 1..%s)" % MAX_ROMAN_NUMERAL
    if int(n) <> n:
        raise NotIntegerError, "non-integers can not be converted"
    return toRomanTable[n]

def fromRoman(s):
    """convert Roman numeral to integer"""
    if not s:
        raise InvalidRomanNumeralError, "Input can not be blank"
    if not fromRomanTable.has_key(s):
        raise InvalidRomanNumeralError, "Invalid Roman numeral: %s" % s
    return fromRomanTable[s]

def toRomanDynamic(n):
```

```
        """convert integer to Roman numeral using dynamic programming"""
    result = ""
    for numeral, integer in romanNumeralMap:
        if n >= integer:
            result = numeral
            n -= integer
            break
    if n > 0:
        result += toRomanTable[n]
    return result

def fillLookupTables():
    """compute all the possible roman numerals"""
    #Save the values in two global tables to convert to and from integers.
    for integer in range(1, MAX_ROMAN_NUMERAL + 1):
        romanNumber = toRomanDynamic(integer)
        toRomanTable.append(romanNumber)
        fromRomanTable[romanNumber] = integer

fillLookupTables()
```

So how fast is it?

**Example 15.18. Output of `romantest9.py` against `roman9.py`**

```
.............
----------------------------------------------------------------------
Ran 13 tests in 0.791s

OK
```

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, because this version will take longer to import (when it fills the lookup tables). But since import is only done once, this is negligible in the long run.

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

- 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 `if __name__` 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/software.htm) for many different languages.

# 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 implementation of basic unit tests in Python. In Chapter 15, *Refactoring*, you saw how unit testing makes large–scale refactoring easier. This chapter will build on those sample programs, but here we will focus more on advanced Python–specific techniques, rather than on unit testing itself.

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 for individual modules, collects them all into one big test suite, and runs them all at once. I actually use this script as part of the build process for this book; I have unit tests for several of the example programs (not just the `roman.py` module featured in Chapter 13, *Unit Testing*), and the first thing my automated build script does is run this program to make sure all my examples still work. If this regression test fails, the build immediately stops. I don't want to release non–working examples any more than you want to download them and sit around scratching your head and yelling at your monitor and wondering why they don't work.

**Example 16.1. `regression.py`**

If you have not already done so, you can download this and other examples (http://diveintopython.org/download/diveintopython–examples–5.4.zip) used in this book.

```
"""Regression testing framework

This module will search for scripts in the same directory named
XYZtest.py.  Each such script should be a test suite that tests a
module through PyUnit.  (As of Python 2.1, PyUnit is included in
the standard library as "unittest".)  This script will aggregate all
found test suites into one big test suite and run them all at once.
"""

import sys, os, re, unittest

def regressionTest():
    path = os.path.abspath(os.path.dirname(sys.argv[0]))
    files = os.listdir(path)
    test = re.compile("test\.py$", re.IGNORECASE)
    files = filter(test.search, files)
    filenameToModuleName = lambda f: os.path.splitext(f)[0]
    moduleNames = map(filenameToModuleName, files)
    modules = map(__import__, moduleNames)
    load = unittest.defaultTestLoader.loadTestsFromModule
    return unittest.TestSuite(map(load, modules))

if __name__ == "__main__":
    unittest.main(defaultTest="regressionTest")
```

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 *moduletest.py*, run them as a single test, and pass or fail them all at once.

**Example 16.2. Sample output of `regression.py`**

```
[you@localhost py]$ python regression.py -v
help should fail with no object ... ok                    ❶
```

```
help should return known result for apihelper ... ok
help should honor collapse argument ... ok
help should honor spacing argument ... ok
buildConnectionString should fail with list input ... ok          ❷
buildConnectionString should fail with string input ... ok
buildConnectionString should fail with tuple input ... ok
buildConnectionString handles empty dictionary ... ok
buildConnectionString returns known result with known input ... ok
fromRoman should only accept uppercase input ... ok               ❸
toRoman should always return uppercase ... ok
fromRoman should fail with blank string ... ok
fromRoman should fail with malformed antecedents ... ok
fromRoman should fail with repeated pairs of numerals ... ok
fromRoman should fail with too many repeated numerals ... ok
fromRoman should give known result with known input ... ok
toRoman should give known result with known input ... ok
fromRoman(toRoman(n))==n for all n ... ok
toRoman should fail with non-integer input ... ok
toRoman should fail with negative input ... ok
toRoman should fail with large input ... ok
toRoman should fail with 0 input ... ok
kgp a ref test ... ok
kgp b ref test ... ok
kgp c ref test ... ok
kgp d ref test ... ok
kgp e ref test ... ok
kgp f ref test ... ok
kgp g ref test ... ok


----------------------------------------------------------------------
Ran 29 tests in 2.799s

OK
```

❶    The first 5 tests are from `apihelpertest.py`, which tests the example script from Chapter 4, *The Power Of Introspection*.

❷    The next 5 tests are from `odbchelpertest.py`, which tests the example script from Chapter 2, *Your First Python Program*.

❸    The rest are from `romantest.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 to it is `sys.argv`. As you saw in Chapter 9, *XML Processing*, this is a list that holds the list of command–line arguments. However, it also holds the name of the running script, exactly as it was called from the command line, and this is enough information to determine its location.

**Example 16.3. `fullpath.py`**

Ifdret ehi7s in 2.799s

```
print 'sys.argv[0] =', sys.argv[0]                    ❶
pathname = os.path.dirname(sys.argv[0])               ❷
print 'path =', pathname
print 'full path =', os.path.abspath(pathname)        ❸
```

❶  Regardless of how you run a script, `sys.argv[0]` will always contain the name of the script, exactly as it appears on the command line. This may or may not include any path information, as you'll see shortly.

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 an empty string.

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.

`os.path.abspath` deserves further explanation. It is very flexible; it can take any kind of pathname.

❶
❷
❸
❹
❺

❶
❷

❸

❹
❺

❶

❷

```
[you@localhost py]$ python fullpath.py                                    ❸
sys.argv[0] = fullpath.py
path =
full path = /home/you/diveintopython/common/py
```

❶   In the first case, `sys.argv[0]` 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 name, and
    `os.path.abspath` simply returns what you give it.

❷   If the script is run by using a partial pathname, `sys.argv[0]` will still contain exactly what appears on the
    command line. `os.path.dirname` will then give you a partial pathname (relative to the current directory),
    and `os.path.abspath` will construct a full pathname from the partial pathname.

❸   If the script is run from the current directory without giving any path, `os.path.dirname` will simply return
    an empty string. Given an empty string, `os.path.abspath` returns the current directory, which is what you
    want, since the script was run from the current directory.

Like the other functions in the `os` and `os.path` modules, `os.path.abspath` is cross–platform. Your results
will look slightly different than my examples if you're running on Windows (which uses backslash as a path
separator) or Mac OS (which uses colons), but they'll still work. That's the whole point of the `os` module.

**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 `regression.py` is located. He suggests this approach instead:

**Example 16.6. Running scripts in the current directory**

```
import sys, os, re, unittest

def regressionTest():
    path = os.getcwd()              ❶
    sys.path.append(path)           ❷
    files = os.listdir(path)        ❸
```

❶   Instead of setting `path` to the directory where the currently running script is located, you set it to the
    current working directory instead. This will be whatever directory you were in before you ran the script,
    which is not necessarily the same as the directory the script is in. (Read that sentence a few times until
    you get it.)

❷   Append this directory to the Python library search path, so that when you dynamically import the unit
    test modules later, Python can find them. You didn't need to do this when `path` was the directory of the
    currently running script, because Python always looks in that directory.

❸   The rest of the function is the same.

This technique will allow you to re–use this `regression.py` script on multiple projects. Just put the script in a
common directory, then change to the project's directory before running it. All of that project's unit tests will be found
and tested, instead of the unit tests in the common directory where `regression.py` 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, which some people feel is more expressive.

Python has a built–in `filter` function which takes two arguments, a function and a list, and returns a list.[7] The
function passed as the first argument to `filter` must itself take one argument, and the list that `filter` returns will
contain all the elements from the list passed to `filter` for which the function passed to `filter` returns true.

Got all that? It's not as difficult as it sounds.

## Example 16.7. Introducing `filter`

```
>>> def odd(n):                          ❶
...     return n % 2
...
>>> li = [1, 2, 3, 5, 9, 10, 256, -3]
>>> filter(odd, li)                      ❷
[1, 3, 5, 9, -3]                         ❸
                                         ❹
```

❶
❷

❸
❹

❶
❷
❸

❶

❷

❸

**Example 16.9. Filtering using list comprehensions instead**

```
files = os.listdir(path)
test = re.compile("test\.py$", re.IGNORECASE)
files = [f for f in files if test.search(f)] ❶
```

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 built-in `map` function. It works much the same way as the `filter` function.

**Example 16.10. Introducing `map`**

```
>>> def double(n):
...     return n*2
...
>>> li = [1, 2, 3, 5, 9, 10, 256, -3]
>>> map(double, li)                          ❶
[2, 4, 6, 10, 18, 20, 512, -6]
>>> [double(n) for n in li]                  ❷
[2, 4, 6, 10, 18, 20, 512, -6]
>>> newlist = []
>>> for n in li:                             ❸
...     newlist.append(double(n))
...
>>> newlist
[2, 4, 6, 10, 18, 20, 512, -6]
```

❶ `map` takes a function and a list[8] and returns a new list by calling the function with each element of the list in order. In this case, the function simply multiplies each element 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 forever.

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**

```
>>> li = [5, 'a', (2, 'b')]
>>> map(double, li)                          ❶
[10, 'aa', (2, 'b', 2, 'b')]
```

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 each type. In this case, the `double` function simply multiplies the given argument by 2, and Python Does The Right Thing depending on the datatype of the argument. For integers, this means actually multiplying it by 2; for strings, it means concatenating the string with itself; for tuples, it means making a new tuple that has all of the elements of the original, then all of the elements of the original again.

All right, enough play time. Let's look at some real code.

**Example 16.12. `map` in `regression.py`**

```
filenameToModuleName = lambda f: os.path.splitext(f)[0]  ❶
moduleNames = map(filenameToModuleName, files)             ❷
```

❶   As you saw in Section 4.7, Using lambda Functions , `lambda` defines an inline function. And as you saw in

❷

❶

❶

❶

❷

❶

❷

❶

❷
❸

```
<module 'os' from 'c:\Python22\lib\os.pyc'>,
<module 're' from 'c:\Python22\lib\re.pyc'>,
<module 'unittest' from 'c:\Python22\lib\unittest.pyc'>]
>>> modules[0].version                          ❹
'2.2.2 (#37, Nov 26 2002, 10:24:37) [MSC 32 bit (Intel)]'
>>> import sys
>>> sys.version
'2.2.2 (#37, Nov 26 2002, 10:24:37) [MSC 32 bit (Intel)]'
```

❶    `moduleNames` 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.

❷    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 (`moduleNames`) and calls the function (`__import__`) over and over, once with each element of the list, builds a list of the return values, and returns the result.

❸    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, where you installed Python, the phase of the moon, etc.)

❹    To drive home the point that these are real modules, let's look at some module attributes. Remember, `modules[0]` *is* the `sys` module, so `modules[0].version` *is* `sys.version`. All the other attributes and methods of these modules are also available. There's nothing magic about the `import` statement, and there's nothing magic about modules. Modules are objects. Everything is an object.

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 within it.

**Example 16.16. The `regressionTest` function**

```
def regressionTest():
    path = os.path.abspath(os.path.dirname(sys.argv[0]))
    files = os.listdir(path)
    test = re.compile("test\.py$", re.IGNORECASE)
    files = filter(test.search, files)
    filenameToModuleName = lambda f: os.path.splitext(f)[0]
    moduleNames = map(filenameToModuleName, files)
    modules = map(__import__, moduleNames)
load = unittest.defaultTestLoader.loadTestsFromModule
return unittest.TestSuite(map(load, modules))
```

Let's look at it line by line, interactively. Assume that the current directory is `c:\diveintopython\py`, which contains the examples that come with this book, including this chapter's script. As you saw in Section 16.2, Finding the path , the script directory will end up in the `path` variable, so let's start hard–code that and go from there.

**Example 16.17. Step 1: Get all the files**

```
>>> import sys, os, re, unittest
>>> path = r'c:\diveintopython\py'
>>> files = os.listdir(path)
>>> files ❶
```

```
['BaseHTMLProcessor.py', 'LICENSE.txt', 'apihelper.py', 'apihelpertest.py',
'argecho.py', 'autosize.py', 'builddialectexamples.py', 'dialect.py',
'fileinfo.py', 'fullpath.py', 'kgptest.py', 'makerealworddoc.py',
'odbchelper.py', 'odbchelpertest.py', 'parsephone.py', 'piglatin.py',
'plural.py', 'pluraltest.py', 'pyfontify.py', 'regression.py', 'roman.py', 'romantest.py',
'uncurly.py', 'unicode2koi8r.py', 'urllister.py', 'kgp', 'plural', 'roman',
'colorize.py']
```

❶ `files` 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 some `.pyc` files in there as well.)

**Example 16.18. Step 2: Filter to find the files you care about**

```
>>> test = re.compile("test\.py$", re.IGNORECASE)          ❶
>>> files = filter(test.search, files)                      ❷
>>> files                                                   ❸
['apihelpertest.py', 'kgptest.py', 'odbchelpertest.py', 'pluraltest.py', 'romantest.py']
```

❶ This regular expression will match any string that ends with `test.py`. Note that you need to escape the period, since a period in a regular expression usually means "match any single character", but you actually want to match a literal period instead.

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 match the regular expression.

❸ And you're left with the list of unit testing scripts, because they were the only ones named `SOMETHINGtest.py`.

**Example 16.19. Step 3: Map filenames to module names**

```
>>> filenameToModuleName = lambda f: os.path.splitext(f)[0]   ❶
>>> filenameToModuleName('romantest.py')                      ❷
'romantest'
>>> filenameToModuleName('odchelpertest.py')
'odbchelpertest'
>>> moduleNames = map(filenameToModuleName, files)            ❸
>>> moduleNames                                               ❹
```

❶

❷

❸
❹

❶
❷

```
<module 'pluraltest' from 'pluraltest.py'>,
<module 'romantest' from 'romantest.py'>]
>>> modules[-1]                                                    ❸
<module 'romantest' from 'romantest.py'>
```

❶ As you saw in Section 16.6, Dynamically importing modules , you can use a combination of `map` and `__import__` to map a list of module names (as strings) into actual modules (which you can call or access like any other module).

❷ `modules` is now a list of modules, fully accessible like any other module.

❸ The last module in the list *is* the `romantest` module, just as if you had said `import romantest`.

**Example 16.21. Step 5: Loading the modules into a test suite**

```
>>> load = unittest.defaultTestLoader.loadTestsFromModule
>>> map(load, modules)                                             ❶
[<unittest.TestSuite tests=[
  <unittest.TestSuite tests=[<apihelpertest.BadInput testMethod=testNoObject>]>,
  <unittest.TestSuite tests=[<apihelpertest.KnownValues testMethod=testApiHelper>]>,
  <unittest.TestSuite tests=[
    <apihelpertest.ParamChecks testMethod=testCollapse>,
    <apihelpertest.ParamChecks testMethod=testSpacing>]>,
    ...
  ]
]
>>> unittest.TestSuite(map(load, modules))                         ❷
```

❶ These are real module objects. Not only can you access them like any other module, instantiate classes and call functions, you can also introspect into the module to figure out which classes and functions it has in the first place. That's what the `loadTestsFromModule` method does: it introspects into each module and returns a `unittest.TestSuite` object for each module. Each `TestSuite` object actually contains a list of `TestSuite` objects, one for each `TestCase` class in your module, and each of those `TestSuite` objects contains a list of tests, one for each test method in your module.

❷ Finally, you wrap the list of `TestSuite` objects into one big test suite. The `unittest` module has no problem traversing this tree of nested test suites within test suites; eventually it gets down to an individual test method and executes it, verifies that it passes or fails, and moves on to the next one.

This introspection process is what the `unittest` module usually does for us. Remember that magic–looking `unittest.main()` function that our individual test modules called to kick the whole thing off? `unittest.main()` actually creates an instance of `unittest.TestProgram`, which in turn creates an instance of a `unittest.defaultTestLoader` and loads it up with the module that called it. (How does it get a reference to the module that called it if you don't give it one? By using the equally–magic `__import__('__main__')` command, which dynamically imports the currently–running module. I could write a book on all the tricks and techniques used in the `unittest` module, but then I'd never finish this one.)

**Example 16.22. Step 6: Telling `unittest` to use your test suite**

```
if __name__ == "__main__":
    unittest.main(defaultTest="regressionTest")                    ❶
```

❶ Instead of letting the `unittest` module do all its magic for us, you've done most of it yourself. You've created a function (`regressionTest`) that imports the modules yourself, calls `unittest.defaultTestLoader` yourself, and wraps it all up in a test suite. Now all you need to do is tell `unittest` that, instead of looking for tests and building a test suite in the usual way, it should just call the `regressionTest` function,

which returns a ready–to–use `TestSuite`.

## 16.8. Summary

The `regression.py` 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.

---

[7] Technically, the second argument to `filter` can be any sequence, including lists, tuples, and custom classes that act like lists by defining the `__getitem__` special method. If possible, `filter` will return the same datatype as you give it, so filtering a list returns a list, but filtering a tuple returns a tuple.

[8] Again, I should point out that `map` can take a list, a tuple, or any object that acts like a sequence. See previous footnote about `filter`.

# 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 2.3. But first, let's talk about how to make plural nouns.

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, and quickly descends into more advanced uses.

English is a schizophrenic language that borrows from a lot of other languages, and the rules for making singular nouns into plural nouns are varied and complex. There are rules, and then there are exceptions to those rules, and then there are exceptions to the exceptions.

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:

1. If a word ends in S, X, or Z, add ES. "Bass" becomes "basses", "fax" becomes "faxes", and "waltz" becomes "waltzes".
2. 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 make a sound that you can hear. So "coach" becomes "coaches" and "rash" becomes "rashes", because you can hear the CH and SH sounds when you say them. But "cheetah" becomes "cheetahs", because the H is silent.
3. 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. So "vacancy" becomes "vacancies", but "day" becomes "days".
4. If all else fails, just add S and hope for the best.

(I know, there are a lot of exceptions. "Man" becomes "men" and "woman" becomes "women", but "human" becomes "humans". "Mouse" becomes "mice" and "louse" becomes "lice", but "house" becomes "houses". "Knife" becomes "knives" and "wife" becomes "wives", but "lowlife" becomes "lowlifes". And don't even get me started on words that are their own plural, like "sheep", "deer", and "haiku".)

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 to add more rules, and you may eventually need to add more languages.

## 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 different combinations of characters, and then do different things to them. This sounds like a job for regular expressions.

**Example 17.1. `plural1.py`**

```
import re

def plural(noun):
    if re.search('[sxz]$', noun):                    ❶
        return re.sub('$', 'es', noun)               ❷
    elif re.search('[^aeioudgkprt]h$', noun):
```

```
        return re.sub('$', 'es', noun)
    elif re.search('[^aeiou]y$', noun):
        return re.sub('y$', 'ies', noun)
    else:
        return noun + '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 exactly one of these characters". So `[sxz]` means "s, or x, or z", but only one of them. The `$` should be familiar; it matches the end of string. So you're checking to see if `noun` ends with s, x, or z.

❷  This `re.sub` function performs regular expression–based string substitutions. Let's look at it in more detail.

### Example 17.2. Introducing `re.sub`

```
>>> import re
>>> re.search('[abc]', 'Mark')        ❶
<_sre.SRE_Match object at 0x001C1FA8>
>>> re.sub('[abc]', 'o', 'Mark')      ❷
'Mork'
>>> re.sub('[abc]', 'o', 'rock')      ❸
'rook'
>>> re.sub('[abc]', 'o', 'caps')      ❹
'oops'
```

❶  Does the string `Mark` contain a, b, or c? Yes, it contains a.

❷  OK, now find a, b, or c, and replace it with o. `Mark` becomes `Mork`.

❸  The same function turns `rock` into `rook`.

❹  You might think this would turn `caps` into `oaps`, but it doesn't. `re.sub` replaces *all* of the matches, not just the first one. So this regular expression turns `caps` into `oops`, because both the c and the a get turned into o.

### Example 17.3. Back to `plural1.py`

```
import re

def plural(noun):
    if re.search('[sxz]$', noun):
        return re.sub('$', 'es', noun)            ❶
    elif re.search('[^aeioudgkprt]h$', noun):     ❷
        return re.sub('$', 'es', noun)            ❸
    elif re.search('[^aeiou]y$', noun):
        return re.sub('y$', 'ies', noun)
    else:
        return noun + 's'
```

❶  Back to the `plural` function. What are you doing? You're replacing the end of string with `es`. In other words, adding `es` to the string. You could accomplish the same thing with string concatenation, for example `noun + 'es'`, but I'm using regular expressions for everything, for consistency, for reasons that will become clear later in the chapter.

❷  Look closely, this is another new variation. The `^` as the first character inside the square brackets means something special: negation. `[^abc]` means "any single character *except* a, b, or c". So `[^aeioudgkprt]` means any character except a, e, i, o, u, d, g, k, p, r, or t. Then that character needs to be followed by h, followed by end of string. You're looking for words that end in H where the H can be heard.

❸  Same pattern here: match words that end in Y, where the character before the Y is *not* a, e, i, o, or u. You're looking for words that end in Y that sounds like I.

**Example 17.4. More on negation regular expressions**

```
>>> import re
>>> re.search('[^aeiou]y$', 'vacancy')  ❶
<_sre.SRE_Match object at 0x001C1FA8>
>>> re.search('[^aeiou]y$', 'boy')      ❷
>>>
>>> re.search('[^aeiou]y$', 'day')
>>>
>>> re.search('[^aeiou]y$', 'pita')     ❸
>>>
```

❶  `vacancy` matches this regular expression, because it ends in `cy`, and `c` is not `a`, `e`, `i`, `o`, or `u`.

❷  `boy` does not match, because it ends in `oy`, and you specifically said that the character before the `y` could not be `o`. `day` does not match, because it ends in `ay`.

❸  `pita` does not match, because it does not end in `y`.

**Example 17.5. More on `re.sub`**

```
>>> re.sub('y$', 'ies', 'vacancy')            ❶
'vacancies'
>>> re.sub('y$', 'ies', 'agency')
'agencies'
>>> re.sub('([^aeiou])y$', r'\1ies', 'vacancy') ❷
'vacancies'
```

❶  This regular expression turns `vacancy` into `vacancies` and `agency` into `agencies`, which is what you wanted. Note that it would also turn `boy` into `boies`, but that will never happen in the function because you did that `re.search` first to find out whether you should do this `re.sub`.

❷  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 another to actually apply it) into a single regular expression. Here's what that would look like. Most of it should look familiar: you're using a remembered group, which you learned in Section 7.6,  Case study: Parsing Phone Numbers , to remember the character before the `y`. Then in the substitution string, you use a new syntax, `\1`, which means "hey, that first group you remembered? put it here". In this case, you remember the `c` before the `y`, and then when you do the substitution, you substitute `c` in place of `c`, and `ies` in place of `y`. (If you have more than one remembered group, you can use `\2` and `\3` and so on.)

Regular expression substitutions are extremely powerful, and the `\1` syntax makes them even more powerful. But combining the entire operation into one regular expression is also much harder to read, and it doesn't directly map to the way you first described the pluralizing rules. You originally laid out rules like "if the word ends in S, X, or Z, then add ES". And if you look at this function, you have two lines of code that say "if the word ends in S, X, or Z, then add ES". It doesn't get much more direct than that.

## 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. Let's temporarily complicate part of the program so you can simplify another part.

**Example 17.6. `plural2.py`**

```
import re

def match_sxz(noun):
```

```python
    return re.search('[sxz]$', noun)
```

❶

❷
❸
❹

❶

❷

❸

❹

```
    if match_y(noun):
        return apply_y(noun)
    if match_default(noun):
        return apply_default(noun)
```

The benefit here is that that `plural` function is now simplified. It takes a list of rules, defined elsewhere, and iterates through them in a generic fashion. Get a match rule; does it match? Then call the apply rule. The rules could be defined anywhere, in any way. The `plural` function doesn't care.

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 previous example, it would require adding an `if` statement to the `plural` function. In this example, it would require adding two functions, `match_foo` and `apply_foo`, and then updating the `rules` list to specify where in the order the new match and apply functions should be called relative to the other rules.

This is really just a stepping stone to the next section. Let's move on.

## 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 define them in the `rules` list and call them through there. Let's streamline the rules definition by anonymizing those functions.


**Example 17.8. `plural3.py`**

```
import re

rules = \
  (
    (
     lambda word: re.search('[sxz]$', word),
     lambda word: re.sub('$', 'es', word)
    ),
    (
     lambda word: re.search('[^aeioudgkprt]h$', word),
     lambda word: re.sub('$', 'es', word)
    ),
    (
     lambda word: re.search('[^aeiou]y$', word),
     lambda word: re.sub('y$', 'ies', word)
    ),
    (
     lambda word: re.search('$', word),
     lambda word: re.sub('$', 's', word)
    )
  )                                                    ❶

def plural(noun):
    for matchesRule, applyRule in rules:               ❷
        if matchesRule(noun):
            return applyRule(noun)
```

❶   This is the same set of rules as you defined in stage 2. The only difference is that instead of defining
     named functions like `match_sxzaptch_sxz`

❷

the first rule, and if it returns a true value, calls the second rule and returns the value. Same as above, word for word. The only difference is that the rule functions were defined inline, anonymously, using lambda functions. But the `plural` function doesn't care how they were defined; it just gets a list of rules and blindly works through them.

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 rule functions inline like this makes it very clear that you have some unnecessary duplication here. You have four pairs of functions, and they all follow the same pattern. The match function is a single call to `re.search`, and the apply function is a single call to `re.sub`. Let's factor out these similarities.

## 17.5. `plural.py`, stage 4

Let's factor out the duplication in the code so that defining new rules can be easier.

**Example 17.9. `plural4.py`**

```
import re

def buildMatchAndApplyFunctions((pattern, search, replace)):
    matchFunction = lambda word: re.search(pattern, word)        ❶
    applyFunction = lambda word: re.sub(search, replace, word)   ❷
    return (matchFunction, applyFunction)                        ❸
```

❶ `buildMatchAndApplyFunctions` is a function that builds other functions dynamically. It takes `pattern`, `search` and `replace` (actually it takes a tuple, but more on that in a minute), and you can build the match function using the `lambda` syntax to be a function that takes one parameter (`word`) and calls `re.search` with the `pattern` that was passed to the `buildMatchAndApplyFunctions` function, and the `word` that was passed to the match function you're building. Whoa.

❷ Building the apply function works the same way. The apply function is a function that takes one parameter, and calls `re.sub` with the `search` and `replace` parameters that were passed to the `buildMatchAndApplyFunctions` function, and the `word` that was passed to the apply function you're building. This technique of using the values of outside parameters within a dynamic function is called *closures*. You're essentially defining constants within the apply function you're building: it takes one parameter (`word`), but it then acts on that plus two other values (`search` and `replace`) which were set when you defined the apply function.

❸ Finally, the `buildMatchAndApplyFunctions` function returns a tuple of two values: the two functions you just created. The constants you defined within those functions (`pattern` within `matchFunction`, and `search` and `replace` within `applyFunction`) stay with those functions, even after you return from `buildMatchAndApplyFunctions`. That's insanely cool.

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. `plural4.py` continued**

```
patterns = \
  (
```

❶
❷

❶

❷

❶

```
def plural(noun, language='en'):                               ❷
    lines = file('rules.%s' % language).readlines()            ❸
    patterns = map(string.split, lines)                        ❹
    rules = map(buildRule, patterns)                           ❺
    for rule in rules:
        result = rule(noun)                                    ❻
        if result: return result
```

❶  You're still using the closures technique here (building a function dynamically that uses variables defined outside the function), but now you've combined the separate match and apply functions into one. (The reason for this change will become clear in the next section.) This will let you accomplish the same thing as having two functions, but you'll need to call it differently, as you'll see in a minute.

❷  Our `plural` function now takes an optional second parameter, `language`, which defaults to `en`.

❸  You use the `language` parameter to construct a filename, then open the file and read the contents into a list. If `language` is `en`, then you'll open the `rules.en` file, read the entire thing, break it up by carriage returns, and return a list. Each line of the file will be one element in the list.

❹  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). Mapping the `string.split` function onto this list will create a new list where each element is a tuple of three strings. So a line like `[sxz]$ $ es` will be broken up into the tuple `('[sxz]$', '$', 'es')`. This means that `patterns` will end up as a list of tuples, just like you hard–coded it in stage 4.

❺  If `patterns` is a list of tuples, then `rules` will be a list of the functions created dynamically by each call to `buildRule`. Calling `buildRule(('[sxz]$', '$', 'es'))` returns a function that takes a single parameter, `word`. When this returned function is called, it will execute `re.search('[sxz]$', word)` and `re.sub('$', 'es', word)`.

❻

```
for applyRule in rules(language):
    result = applyRule(noun)
    if result: return result
```

❶

❷
❸

❹

❺

❻

❶

❷

❸
❹

❺

❻

❶

❷
❸

❶ 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 up slowly at first, then more and more rapidly. To start the sequence, you need two variables: a starts at 0, and b starts at 1.

❷ a is the current number in the sequence, so yield it.

❸ b is the next number in the sequence, so assign that to a, but also calculate the next value (a+b) and assign that to b for later use. Note that this happens in parallel; if a is 3 and b is 5, then a, b = b, a+b will set a to 5 (the previous value of b) and b to 8 (the sum of the previous values of a and b).

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, it works well with for loops.

### Example 17.20. Generators in `for` loops

```
>>> for n in fibonacci(1000):     ❶
...     print n,                  ❷
0 1 1 2 3 5 8 13 21 34 55 89 144 233 377 610 987
```

❶ You can use a generator like fibonacci in a for loop directly. The for loop will create the generator object and successively call the next() method to get values to assign to the for loop index variable (n).

❷ Each time through the for loop, n gets a new value from the yield statement in fibonacci, and all you do is print it out. Once fibonacci runs out of numbers (a gets bigger than max, which in this case is 1000), then the for loop exits gracefully.

OK, let's go back to the plural function and see how you're using this.

### Example 17.21. Generators that generate dynamic functions

```
def rules(language):
    for line in file('rules.%s' % language):                          ❶
        pattern, search, replace = line.split()                       ❷
        yield lambda word: re.search(pattern, word) and re.sub(search, replace, word)   ❸

def plural(noun, language='en'):
    for applyRule in rules(language):     ❹
        result = applyRule(noun)
        if result: return result
```

❶ `for line in file(...)` is a common idiom for reading lines from a file, one line at a time. It works because *file actually returns a generator* whose next() method returns the next line of the file. That is so insanely cool, I wet myself just thinking about it.

❷ No magic here. Remember that the lines of the rules file have three values separated by whitespace, so line.split() returns a tuple of 3 values, and you assign those values to 3 local variables.

❸ *And then you yield.* What do you yield? A function, built dynamically with lambda, that is actually a closure (it uses the local variables pattern, search, and replace as constants). In other words, rules is a generator that spits out rule functions.

❹ 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

# 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 optimization. But don't worry about it *too* much.

## 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 application, how much time is going to be spent running that code, compared to the time spent

```
    # 4. remove all "9"s
    digits3 = re.sub('9', '', digits2)

    # 5. pad end with "0"s to 4 characters
    while len(digits3) < 4:
        digits3 += "0"

    # 6. return first 4 characters
    return digits3[:4]

if __name__ == '__main__':
    from timeit import Timer
    names = ('Woo', 'Pilgrim', 'Flingjingwaller')
    for name in names:
        statement = "soundex('%s')" % name
        t = Timer(statement, "from __main__ import soundex")
        print name.ljust(15), soundex(name), min(t.repeat())
```

**Further Reading on Soundex**

- Soundexing and Genealogy (http://www.avotaynu.com/soundex.html) gives a chronology of the evolution of the Soundex and its regional variations.

## 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 running in the background? Are you sure? Every modern computer has background processes running, some all the time, some intermittently. Cron jobs fire off at consistent intervals; background services occasionally "wake up" to do useful things like check for new mail, connect to instant messaging servers, check for application updates, scan for viruses, check whether a disk has been inserted into your CD drive in the last 100 nanoseconds, and so on. Before you start your timing tests, turn everything off and disconnect from the network. Then turn off all the things you forgot to turn off the first time, then turn off the service that's incessantly checking whether the network has come back yet, then ...

And then there's the matter of the variations introduced by the timing framework itself. Does the Python interpreter cache method name lookups? Does it cache code block compilations? Regular expressions? Will your code have side effects if run more than once? Don't forget that you're dealing with small fractions of a second, so small mistakes in your timing framework will irreparably skew your results.

The Python community has a saying: "Python comes with batteries included." Don't write your own timing framework. Python 2.3 comes with a perfectly good one called `timeit`.

**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−5.4.zip) used in this book.

```
>>> import timeit
>>> t = timeit.Timer("soundex.soundex('Pilgrim')",
...     "import soundex")          ❶
>>> t.timeit()                     ❷
```

```
8.21683733547
>>> t.repeat(3, 2000000)      ❸
[16.48319309109, 16.46128984923, 16.44203948912]
```

❶    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; in this case, you are timing a call to the Soundex function within the `soundex` with an argument of `'Pilgrim'`. The second argument to the `Timer` class is the import statement that sets up the environment for the statement. Internally, `timeit` sets up an isolated virtual environment, manually executes the setup statement (importing the `soundex` module), then manually compiles and executes the timed statement (calling the Soundex function).

❷    Once you have the `Timer` object, the easiest thing to do is call `timeit()`, which calls your function 1

❸

How does `soundex1a.py` perform? For convenience, the `__main__` section of the script contains this code that calls the `timeit` module, sets up a timing test with three different names, tests each name three times, and displays the minimum time for each:

```
if __name__ == '__main__':
    from timeit import Timer
    names = ('Woo', 'Pilgrim', 'Flingjingwaller')
    for name in names:
        statement = "soundex('%s')" % name
        t = Timer(statement, "from __main__ import soundex")
        print name.ljust(15), soundex(name), min(t.repeat())
```

So how does `soundex1a.py` perform with this regular expression?

```
C:\samples\soundex\stage1>python soundex1a.py
Woo             W000 19.3356647283
Pilgrim         P426 24.0772053431
Flingjingwaller F452 35.0463220884
```

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 that gap (make the function take less relative time for longer input), but the nature of the algorithm dictates that it will never run in constant time.

The other thing to keep in mind is that we are testing a representative sample of names. `Woo`

But is this the wrong path? The logic here is simple: the input `source` needs to be non–empty, and it needs to be composed entirely of letters. Wouldn't it be faster to write a loop checking each character, and do away with regular expressions altogether?

Here is `soundex/stage1/soundex1d.py`:

```
if not source:
    return "0000"
for c in source:
    if not ('A' <= c <= 'Z') and not ('a' <= c <= 'z'):
        return "0000"
```

It turns out that this technique in `soundex1d.py` is *not* faster than using a compiled regular expression (although it is faster than using a non–compiled regular expression):

```
C:\samples\soundex\stage1>python soundex1d.py
Woo             W000 15.4065058548
Pilgrim         P426 22.2753567842
Flingjingwaller F452 37.5845122774
```

Why isn't `soundex1d.py` faster? The answer lies in the interpreted nature of Python. The regular expression engine is written in C, and compiled to run natively on your computer. On the other hand, this loop is written in Python, and runs through the Python interpreter. Even though the loop is relatively simple, it's not simple enough to make up for the overhead of being interpreted. Regular expressions are never the right answer... except when they are.

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 book. The method is called `isalpha()`, and it checks whether a string contains only letters.

This is `soundex/stage1/soundex1e.py`:

```
if (not source) and (not source.isalpha()):
    return "0000"
```

How much did we gain by using this specific method in `soundex1e.py`? Quite a bit.

```
C:\samples\soundex\stage1>python soundex1e.py
Woo             W000 13.5069504644
Pilgrim         P426 18.2199394057
Flingjingwaller F452 28.9975225902
```

**Example 18.3. Best Result So Far: `soundex/stage1/soundex1e.py`**

```
import string, re

charToSoundex = {"A": "9",
                 "B": "1",
                 "C": "2",
                 "D": "3",
                 "E": "9",
                 "F": "1",
                 "G": "2",
                 "H": "9",
                 "I": "9",
                 "J": "2",
                 "K": "2",
                 "L": "4",
                 "M": "5",
```

```
                    "N": "5",
                    "O": "9",
                    "P": "1",
                    "Q": "2",
                    "R": "6",
                    "S": "2",
                    "T": "3",
                    "U": "9",
                    "V": "1",
                    "W": "9",
                    "X": "2",
                    "Y": "9",
                    "Z": "2"}

def soundex(source):
    if (not source) and (not source.isalpha()):
        return "0000"
    source = source[0].upper() + source[1:]
    digits = source[0]
    for s in source[1:]:
        s = s.upper()
        digits += charToSoundex[s]
    digits2 = digits[0]
    for d in digits[1:]:
        if digits2[-1] != d:
            digits2 += d
    digits3 = re.sub('9', '', digits2)
    while len(digits3) < 4:
        digits3 += "0"
    return digits3[:4]

if __name__ == '__main__':
    from timeit import Timer
    names = ('Woo', 'Pilgrim', 'Flingjingwaller')
    for name in names:
        statement = "soundex('%s')" % name
        t = Timer(statement, "from __main__ import soundex")
        print name.ljust(15), soundex(name), min(t.repeat())
```

# 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 on each character. This is what we have in `soundex/stage1/soundex1c.py` (the current best result so far):

```
charToSoundex = {"A": "9",
                 "B": "1",
                 "C": "2",
                 "D": "3",
                 "E": "9",
                 "F": "1",
                 "G": "2",
                 "H": "9",
                 "I": "9",
                 "J": "2",
                 "K": "2",
                 "L": "4",
                 "M": "5",
```

```
                "N": "5",
                "O": "9",
                "P": "1",
                "Q": "2",
                "R": "6",
                "S": "2",
                "T": "3",
                "U": "9",
                "V": "1",
                "W": "9",
                "X": "2",
                "Y": "9",
                "Z": "2"}

def soundex(source):
    # ... input check omitted for brevity ...
    source = source[0].upper() + source[1:]
    digits = source[0]
    for s in source[1:]:
        s = s.upper()
        digits += charToSoundex[s]
```

You timed `soundex1c.py` already; this is how it performs:

```
C:\samples\soundex\stage1>python soundex1c.py
Woo               W000 14.5341678901
Pilgrim           P426 19.2650071448
Flingjingwaller   F452 30.1003563302
```

This code is straightforward, but is it the best solution? Calling `upper()` on each individual character seems inefficient; it would probably be better to call `upper()` once on the entire string.

Then there's the matter of incrementally building the `digits` string. Incrementally building strings like this is horribly inefficient; internally, the Python interpreter needs to create a new string each time through the loop, then discard the old one.

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 the string method `join()`.

Here is `soundex/stage2/soundex2a.py`, which converts letters to digits by using ¦ and `lambda`:

```
def soundex(source):
    # ...
    source = source.upper()
    digits = source[0] + "".join(map(lambda c: charToSoundex[c], source[1:]))
```

Surprisingly, `soundex2a.py`

```
    source = source.upper()
    digits = source[0] + "".join([charToSoundex[c] for c in source[1:]])
```

Using a list comprehension in `soundex2b.py` is faster than using | and `lambda` in `soundex2a.py`, but still not faster than the original code (incrementally building a string in `soundex1c.py`):

```
C:\samples\soundex\stage2>python soundex2b.py
Woo             W000 13.4221324219
Pilgrim         P426 16.4901234654
Flingjingwaller F452 25.8186157738
```

It's time for a radically different approach. Dictionary lookups are a general purpose tool. Dictionary keys can be any length string (or many other data types), but in this case we are only dealing with single–character keys *and* single–character values. It turns out that Python has a specialized function for handling exactly this situation: the `string.maketrans` function.

This is `soundex/stage2/soundex2c.py`:

```
allChar = string.uppercase + string.lowercase
charToSoundex = string.maketrans(allChar, "91239129922455912623919292" * 2)
def soundex(source):
    # ...
    digits = source[0].upper() + source[1:].translate(charToSoundex)
```

What the heck is going on here? `string.maketrans` creates a translation matrix between two strings: the first argument and the second argument. In this case, the first argument is the string `ABCDEFGHIJKLMNOPQRSTUVWXYZabcdefghijklmnopqrstuvwxyz`, and the second argument is the string `9123912992245591262391929291239129922455912623919292`. See the pattern? It's the same conversion pattern we were setting up longhand with a dictionary. A maps to 9, B maps to 1, C maps to 2, and so forth. But it's not a dictionary; it's a specialized data structure that you can access using the string method `translate`, which translates each character into the corresponding digit, according to the matrix defined by `string.maketrans`.

`timeit` shows that `soundex2c.py` is significantly faster than defining a dictionary and looping through the input and building the output incrementally:

```
C:\samples\soundex\stage2>python soundex2c.py
Woo             W000 11.437645008
Pilgrim         P426 13.2825062962
Flingjingwaller F452 18.5570110168
```

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: `soundex/stage2/soundex2c.py`**

```
import string, re

allChar = string.uppercase + string.lowercase
charToSoundex = string.maketrans(allChar, "91239129922455912623919292" * 2)
isOnlyChars = re.compile('^[A-Za-z]+$').search

def soundex(source):
    if not isOnlyChars(source):
        return "0000"
    digits = source[0].upper() + source[1:].translate(charToSoundex)
```

```
        digits2 = digits[0]
        for d in digits[1:]:
            if digits2[-1] != d:
                digits2 += d
        digits3 = re.sub('9', '', digits2)
        while len(digits3) < 4:
            digits3 += "0"
        return digits3[:4]

if __name__ == '__main__':
    from timeit import Timer
    names = ('Woo', 'Pilgrim', 'Flingjingwaller')
    for name in names:
        statement = "soundex('%s')" % name
        t = Timer(statement, "from __main__ import soundex")
        print name.ljust(15), soundex(name), min(t.repeat())
```

# 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 `soundex/stage2/soundex2c.py`:

```
        digits2 = digits[0]
        for d in digits[1:]:
            if digits2[-1] != d:
                digits2 += d
```

Here are the performance results for `soundex2c.py`:

```
C:\samples\soundex\stage2>python soundex2c.py
Woo             W000 12.6070768771
Pilgrim         P426 14.4033353401
Flingjingwaller F452 19.7774882003
```

The first thing to consider is whether it's efficient to check `digits[-1]` each time through the loop. Are list indexes expensive? Would we be better off maintaining the last digit in a separate variable, and checking that instead?

To answer this question, here is `soundex/stage3/soundex3a.py`:

```
        digits2 = ''
        last_digit = ''
        for d in digits:
            if d != last_digit:
                digits2 += d
                last_digit = d
```

`soundex3a.py` does not run any faster than `soundex2c.py`, and may even be slightly slower (although it's not enough of a difference to say for sure):

```
C:\samples\soundex\stage3>python soundex3a.py
Woo             W000 11.5346048171
Pilgrim         P426 13.3950636184
Flingjingwaller F452 18.6108927252
```

Why isn't `soundex3a.py` faster? It turns out that list indexes in Python are extremely efficient. Repeatedly accessing `digits2[-1]` is no problem at all. On the other hand, manually maintaining the last seen digit in a separate variable means we have *two* variable assignments for each digit we're storing, which wipes out any small

gains we might have gotten from eliminating the list lookup.

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 to iterate through the list. The problem is, the code needs access to the previous character in the list, and that's not easy to do with a straightforward 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 the list and pull out each character that is different from the previous character. That will give you a list of characters, and you can use the string method `join()` to reconstruct a string from that.

Here is `soundex/stage3/soundex3b.py`:

```
digits2 = "".join([digits[i] for i in range(len(digits))
                   if i == 0 or digits[i-1] != digits[i]])
```

Is this faster? In a word, no.

```
C:\samples\soundex\stage3>python soundex3b.py
Woo               W000 14.2245271396
Pilgrim           P426 17.8337165757
Flingjingwaller   F452 25.9954005327
```

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: `list('abc')` returns `['a', 'b', 'c']`. Furthermore, lists can be *modified in place*its–centr)ach cracring ourcng metho, whyasy tmov –13.2 Td(list, e, t". ss the ca strce)ng comman'abc?–13.2 Td(Let's H

```
allChar = string.uppercase + string.lowercase
charToSoundex = string.maketrans(allChar, "91239129922455912623919292" * 2)
isOnlyChars = re.compile('^[A-Za-z]+$').search

def soundex(source):
    if not isOnlyChars(source):
        return "0000"
    digits = source[0].upper() + source[1:].translate(charToSoundex)
    digits2 = digits[0]
    for d in digits[1:]:
        if digits2[-1] != d:
            digits2 += d
    digits3 = re.sub('9', '', digits2)
    while len(digits3) < 4:
        digits3 += "0"
    return digits3[:4]

if __name__ == '__main__':
    from timeit import Timer
    names = ('Woo', 'Pilgrim', 'Flingjingwaller')
    for name in names:
        statement = "soundex('%s')" % name
        t = Timer(statement, "from __main__ import soundex")
        print name.ljust(15), soundex(name), min(t.repeat())
```

## 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 `soundex/stage2/soundex2c.py`:

```
    digits3 = re.sub('9', '', digits2)
    while len(digits3) < 4:
        digits3 += "0"
    return digits3[:4]
```

These are the results for `soundex2c.py`:

```
C:\samples\soundex\stage2>python soundex2c.py
Woo             W000 12.6070768771
Pilgrim         P426 14.4033353401
Flingjingwaller F452 19.7774882003
```

The first thing to consider is replacing that regular expression with a loop. This code is from `soundex/stage4/soundex4a.py`:

```
    digits3 = ''
    for d in digits2:
        if d != '9':
            digits3 += d
```

Is `soundex4a.py` faster? Yes it is:

```
C:\samples\soundex\stage4>python soundex4a.py
Woo             W000 6.62865531792
Pilgrim         P426 9.02247576158
Flingjingwaller F452 13.6328416042
```

# 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 compiled in C and runs natively on your computer; your loop is written in Python and runs through the Python interpreter.
- If you need to choose between regular expressions and string methods, choose string methods. Both are compiled in C, so choose the simpler one.
- General–purpose dictionary lookups are fast, but specialtiy functions such as `string.maketrans` and string methods such as `isalpha()` are faster. If Python has a custom–tailored function for you, use it.
- 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 over 1 million function calls. Great. Now think: over the course of those million function calls, how many seconds will your surrounding application wait for a database connection? Or wait for disk I/O? Or wait for user input? Don't spend too much time over–optimizing one algorithm, or you'll ignore obvious improvements somewhere else. Develop an instinct for the sort of code that Python runs well, correct obvious blunders if you find them, and leave the rest alone.

# 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 `doc string` conventions.
    - *Python Style Guide* (http://www.python.org/doc/essays/styleguide.html) discusses how to write a good `doc string`.
    - *Python Tutorial* (http://www.python.org/doc/current/tut/tut.html) discusses conventions for spacing in `doc strings` (http://www.python.org/doc/current/tut/node6.html#SECTION006750000000000000000).
- 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 (http://www.python.org/doc/current/ref/objects.html), because some people are pedantic and like to discuss this sort of thing at great length.
    - 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 various indentation errors (http://www.python.org/doc/current/ref/indentation.html).
    - *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 use dictionaries to model sparse matrices (http://www.ibiblio.org/obp/thinkCSpy/chap10.htm).
    - Python Knowledge Base (http://www.faqts.com/knowledge–base/index.phtml/fid/199/) has a lot of example code using dictionaries (http://www.faqts.com/knowledge–base/index.phtml/fid/541).
    - Python Cookbook (http://www.activestate.com/ASPN/Python/Cookbook/) discusses how to sort the values of a dictionary by key (http://www.activestate.com/ASPN/Python/Cookbook/Recipe/52306).
    - *Python Library Reference* (http://www.python.org/doc/current/lib/) summarizes all the dictionary methods (http://www.python.org/doc/current/lib/typesmapping.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 point about passing lists as function arguments (http://www.ibiblio.org/obp/thinkCSpy/chap08.htm).

- ♦ *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/node7.html#SECTION007110000000000000000).
- ♦ Python Knowledge Base (http://www.faqts.com/knowledge−base/index.phtml/fid/199/) answers common questions about lists (http://www.faqts.com/knowledge−base/index.phtml/fid/534) and has a lot of example code using lists (http://www.faqts.com/knowledge−base/index.phtml/fid/540).
- ♦ *Python Library Reference* (http://www.python.org/doc/current/lib/) summarizes all the list methods (http://www.python.org/doc/current/lib/typesseq−mutable.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 concatenate tuples (http://www.ibiblio.org/obp/thinkCSpy/chap10.htm).
  - ♦ Python Knowledge Base (http://www.faqts.com/knowledge−base/index.phtml/fid/199/) shows how to sort a tuple (http://www.faqts.com/knowledge−base/view.phtml/aid/4553/fid/587).
  - ♦ *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/node7.html#SECTION007300000000000000000).

- 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 (http://www.python.org/doc/current/ref/implicit−joining.html) and when you need to use it (http://www.python.org/doc/current/ref/explicit−joining.html).
  - ♦ *How to Think Like a Computer Scientist* (http://www.ibiblio.org/obp/thinkCSpy/) shows how to use multi−variable assignment to swap the values of two variables (http://www.ibiblio.org/obp/thinkCSpy/chap09.htm).

- 3.5. Formatting Strings

  - ♦ *Python Library Reference* (http://www.python.org/doc/current/lib/) summarizes all the string formatting format characters (http://www.python.org/doc/current/lib/typesseq−strings.html).
  - ♦ *Effective AWK Programming* (http://www−gnats.gnu.org:8080/cgi−bin/info2www?(gawk)Top) discusses all the format characters (http://www−gnats.gnu.org:8080/cgi−bin/info2www?(gawk)Control+Letters) and advanced string formatting techniques like specifying width, precision, and zero−padding (http://www−gnats.gnu.org:8080/cgi−bin/info2www?(gawk)Format+Modifiers).

- 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 (http://www.python.org/doc/current/tut/node7.html#SECTION007130000000000000000).
  - ♦ *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/node7.html#SECTION007140000000000000000).

3.4://www.faq:1.033·

(http://www.python.org/cgi−bin/faqw.py?query=4.96&querytype=simple&casefold=yes&req=search) instead of a list method.

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 (http://www.python.org/doc/current/tut/node6.html#SECTION006710000000000000000), which matters when the default value is a list or an expression with side effects.
- 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−funcs.html) and all the built−in exceptions (http://www.python.org/doc/current/lib/module−exceptions.html).
- 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/node7.html#SECTION007130000000000000000).
- 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/52310).
- 4.7.1. Real−World lambda Functions

    ◆ Python Knowledge Base (http://www.faqts.com/knowledge−base/index.phtml/fid/199/) discusses using `lambda` to call functions indirectly (http://www.faqts.com/knowledge−base/view.phtml/aid/6081/fid/241).
    ◆ *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/node6.html#SECTION006740000000000000000). (PEP 227 (http://python.sourceforge.net/peps/pep−0227.html) explains how this will change in future versions of Python.)
    ◆ *The Whole Python FAQ* (http://www.python.org/doc/FAQ.html) has examples of obfuscated one−liners using `lambda` (http://www.python.org/cgi−bin/faqw.py?query=4.15&querytype=simple&casefold=yes&req=search).

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 `import` *module vs.* `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/node8.html#SECTION008410000000000000000).
- 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/tutclass.htm).

- ♦ *How to Think Like a Computer Scientist* (http://www.ibiblio.org/obp/thinkCSpy/) shows how to use classes to model compound datatypes (http://www.ibiblio.org/obp/thinkCSpy/chap12.htm).
  - ♦ *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/node11.html).
  - ♦ Python Knowledge Base (http://www.faqts.com/knowledge–base/index.phtml/fid/199/) answers common questions about classes (http://www.faqts.com/knowledge–base/index.phtml/fid/242).
- 5.4.1. Garbage Collection

  - ♦ *Python Library Reference* (http://www.python.org/doc/current/lib/) summarizes built–in attributes like __class__ (http://www.python.org/doc/current/lib/specialattrs.html).
  - ♦ *Python Library Reference* (http://www.python.org/doc/current/lib/) documents the gc module (http://www.python.org/doc/current/lib/module–gc.html), which gives you low–level control over Python's garbage collection.
- 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/module–UserDict.html) and the copy module (http://www.python.org/doc/current/lib/module–copy.html).
- 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/specialnames.html).
- 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/node11.html#SECTION0011600000000000000000).

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 handling multiple exceptions at once (http://www.python.org/doc/current/tut/node10.html#SECTION0010400000000000000000).
  - ♦ *Python Library Reference* (http://www.python.org/doc/current/lib/) summarizes all the built–in exceptions (http://www.python.org/doc/current/lib/module–exceptions.html).
  - ♦ *Python Library Reference* (http://www.python.org/doc/current/lib/) documents the getpass (http://www.python.org/doc/current/lib/module–getpass.html) module.
  - ♦ *Python Library Reference* (http://www.python.org/doc/current/lib/) documents the traceback module (http://www.python.org/doc/current/lib/module–traceback.html), which provides low–level access to exception attributes after an exception is raised.
  - ♦ *Python Reference Manual* (http://www.python.org/doc/current/ref/) discusses the inner workings of the try...except block (http://www.python.org/doc/current/ref/try.html).
- 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 a file one line at a time into a list (http://www.python.org/doc/current/tut/node9.html#SECTION0092100000000000000000).
  - ♦ eff–bot (http://www.effbot.org/guides/) discusses efficiency and performance of various ways of reading a file (http://www.effbot.org/guides/readline–performance.htm).
  - ♦ Python Knowledge Base (http://www.faqts.com/knowledge–base/index.phtml/fid/199/) answers

common questions about files (http://www.faqts.com/knowledge−base/index.phtml/fid/552).

♦ *Python Library Reference* (http://www.python.org/doc/current/lib/) summarizes all the file object methods (http://www.python.org/doc/current/lib/bltin−file−objects.html).

- 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/node6.html#SECTION006710000000000000000).
    ♦ *Python Library Reference* (http://www.python.org/doc/current/lib/) documents the `sys` (http://www.python.org/doc/current/lib/module−sys.html) module.

- 6.5. Working with Directories

    ♦ Python Knowledge Base (http://www.faqts.com/knowledge−base/index.phtml/fid/199/) answers questions about the `os` module (http://www.faqts.com/knowledge−base/index.phtml/fid/240).
    ♦ *Python Library Reference* (http://www.python.org/doc/current/lib/) documents the `os` (http://www.python.org/doc/current/lib/module−os.html) module and the `os.path` (http://www.python.org/doc/current/lib/module−os.path.html) 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 in Python.
    ♦ *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#entities).
    ♦ *Python Library Reference* (http://www.python.org/doc/current/lib/) confirms your suspicions that the `htmlentitydefs` module (http://www.python.org/doc/current/lib/module−htmlentitydefs.html) is exactly what it sounds like.

- 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://rinkworks.com/dialect/). Unfortunately, source code does not appear to be available.

## 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 (http://www.unicode.org/standard/principles.html).
    ♦ Unicode Tutorial (http://www.reportlab.com/i18n/python_unicode_tutorial.html) has some more examples of how to use Python's unicode functions, including how to force Python to coerce unicode into ASCII even when it doesn't really want to.
    ♦ PEP 263 (http://www.python.org/peps/pep−0263.html) goes into more detail about how and when to define a character encoding in your `.py` files.

♦ XProgramming.com (http://www.xprogramming.com/) has links to download unit testing frameworks (http://www.xprogramming.com/software.htm) for many different languages.

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://www.google.com/search?q=generators+cookbook+site:aspn.activestate.com).

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 and its regional variations.

# 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

You can document a Python function by giving it a `doc string`.

- 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` or `end`, and no curly braces to mark where the function code starts and stops. The only delimiter is a colon (`:`) and the indentation of the code itself.

- 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 that uses the `if __name__` trick.

Chapter 3. Native Datatypes

    3.1. Introducing Dictionaries

for.

    Python has two ways of importing modules. Both are useful, and you should know when to
    use each. One way, `import module`, you've already seen in Section 2.4, Everything Is an
    Object . The other way accomplishes the same thing, but it has subtle and important
    differences.

    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.

    Instantiating classes in Python is straightforward. To instantiate a class, simply call the class
    as if it were a function, passing the arguments that the `__init__` method defines. The return
    value will be the newly created object.

    As you've seen, `FileInfo` is a class that acts like a dictionary. To explore this further, let's
    look at the `UserDict` class in the `UserDict` module, which is the ancestor of the
    `FileInfo` class. This is nothing special; the class is written in Python and stored in a `.py`
    file, just like any other Python code. In particular, it's stored in the `lib` directory in your
    Python installation.

    In addition to normal class methods, there are a number of special methods that Python
    classes can define. Instead of being called directly by your code (like normal methods),
    special methods are called for you by Python in particular circumstances or when specific
    syntax is used.

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

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

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

    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 `getattr` to create a proxy to a remote web service.

Chapter 6. Exceptions and File Handling

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

some people find more readable. First look at the method we already used in the previous example.

- 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 figure out what one does, that's no guarantee that you'll be able to understand it six months later. What you really need is inline documentation.

- 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 than that. When a regular expression *does* match, you can pick out specific pieces of it. You can find out what matched where.

- 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 them now, believe me, you ain't seen nothing yet.

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|images|links] in my HTML document?" "How do I parse/translate/munge the text of my HTML document but leave the tags alone?" "How can I add/remove/quote attributes of all my HTML tags at once?" This chapter will answer all of these questions.

- 8.2. Introducing sgmllib.py

  HTML processing is broken into three steps: breaking down the HTML into its constituent pieces, fiddling with the pieces, and reconstructing the pieces into HTML again. The first step is done by `sgmllib.py`, a part of the standard Python library.

- 8.3. Extracting data from HTML documents

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

  8.4. Introducing BaseHTMLProcessor.py

  `SGMLParser` 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 don't do anything. `SGMLParser` is an HTML *consumer*: it takes HTML and breaks it down into small, structured pieces. As you saw in the previous section, you can subclass `SGMLParser`

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

A common question on comp.lang.python (http://groups.google.com/groups?group=comp.lang.python) is "I have a bunch of HTML documents with unquoted attribute values, and I want to properly quote them all. How can I do this?"[4] (This is generally precipitated by a project manager who has found the HTML–is–a–standard religion joining a large project and proclaiming that all pages must validate against an HTML validator. Unquoted attribute values are a common violation of the HTML standard.) Whatever the reason, unquoted attribute values are easy to fix by feeding HTML through `BaseHTMLProcessor`.

`Dialectizer` is a simple (and silly) descendant of `BaseHTMLProcessor`. It runs blocks of text through a series of substitutions, but it makes sure that anything within a `<pre>...</pre>` block passes through unaltered.

Traversing XML documents by stepping through each node can be tedious. If you're looking for something in particular, buried deep within your XML document, there is a shortcut you can use to find it quickly: `getElementsByTagName`.

- 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.

- 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 make the program more flexible: using streams for input processing, using `getattr` for method dispatching, and using command–line flags to allow users to reconfigure the program without changing the code.


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 the rest of you.

- 10.3. Caching node lookups

  `kgp.py` 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 of the input documents to build a cache of nodes.

- 10.4. Finding direct children of a node

  Another useful techique when parsing XML documents is finding all the direct child elements of a particular element. For instance, in the grammar files, a `ref` element can have several `p` elements, each of which can contain many things, including other `p` elements. You want to find just the `p` elements that are children of the `ref`, not `p` elements that are children of other `p` elements.

- 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 made up of various types of nodes, each represented by a Python object. The root level of the document itself is represented by a `Document` object. The `Document` then contains one or more `Element` objects (for actual XML tags), each of which may contain other `Element` objects, `Text` objects (for bits of text), or `Comment` objects (for embedded comments). Python makes it easy to write a dispatcher to separate the logic for each node type.

- 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–style flags to specify various options. None of this is XML–specific, but this script makes good use of command–line processing, so it seemed like a good time to mention it.

- 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 `minidom` takes an XML file and parses it into Python objects, providing for random access to arbitrary elements. Furthermore, this chapter shows how Python can be used to create a "real" standalone command–line script, complete with command–line flags, command–line arguments, error handling, even the ability to take input from the piped result of a previous program.

## 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 from a URL, but let's dive into the more general topic of HTTP web services.

- 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 over and over again, every hour, to get the latest news from the site that's offering the news feed. Let's do it the quick–and–dirty way first, and then see how you can do better.

- 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 throughout the chapter, as you add more and more features.

  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`

how they all fit together.

- 11.10. Summary

  The `openanything.py` 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 results? Now you can. Here is a program to search Google from Python.

- 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 simple functions for demonstration purposes.

- 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 `SOAPProxy` class proxies local method calls and transparently turns then into invocations of remote SOAP methods. As you've seen, this is a lot of work, and `SOAPProxy` does it quickly and transparently. What it doesn't do is provide any means of method introspection.

- 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 history entirely, since it bores me to tears. There were other standards that tried to do similar things, but WSDL won, so let's learn how to use it.

- 12.7. Searching Google

  Let's finally turn to the sample code that you saw that the beginning of this chapter, which does something more useful and exciting than get the current temperature.

- 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. This chapter has touched on some of the simpler use cases.

## 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 Python under your belt, you're going to step back and look at the steps that happen *before* the code gets written.

- 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: you're going to write a test suite that puts these functions through their paces and makes sure that they behave the way you want them to. You read that right: you're going to write code that tests code that you haven't written yet.

- 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 be in `roman.py`. It is not immediately obvious how it all fits together; none of these classes or methods reference any of the others. There are good reasons for this, as you'll see shortly.

- 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.

- 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 sort of failure; they must fail in the way you expect.

- 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 A to B and the other converts B to A. In these cases, it is useful to create a "sanity check" to make sure that you can convert A to B and back to A without losing precision, incurring rounding errors, or triggering any other sort of bug.

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 going to do this in stages, so you can see all the unit tests fail, then watch them pass one by one as you fill in the gaps in `roman.py`.

- 14.2. roman.py, stage 2

  Now that you have the framework of the `roman` module laid out, it's time to start writing code and passing test cases.

- 14.3. roman.py, stage 3

  Now that `toRoman` behaves correctly with good input (integers from `1` to `3999`), it's time to make it behave correctly with bad input (everything else).

- 14.4. roman.py, stage 4

  Now that `toRoman` is done, it's time to start coding `fromRoman`. Thanks to the rich data structure that maps individual Roman numerals to integer values, this is no more difficult than the `toRoman` function.

In Chapter 13, *Unit Testing*, you learned about the philosophy of unit testing. In Chapter 14, *Test–First Programming*, you stepped through the implementation of basic unit tests in Python. In Chapter 15, *Refactoring*, you saw how unit testing makes large–scale refactoring easier. This chapter will build on those sample programs, but here we will focus more on advanced Python–specific techniques, rather than on unit testing itself.

- 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.

- 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 feel is more expressive.

- 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 built–in `map` function. It works much the same way as the `filter` function.

- 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 that's a perfectly valid question. Mostly, it's a matter of perspective. Using `map` and `filter` forces you to center your thinking around your data.

- 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 selected modules within it.

- 16.8. Summary

  The `regression.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. Generators are new in Python 2.3. But first, let's talk about how to make plural nouns.

- 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 fi6.4l1iffre t tompbintion  of characters. and gher di deiffre t Tj 0 –13.2 Td(ehing, to

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` list and call them through there. Let's streamline the rules definition by anonymizing those functions.

Let's factor out the duplication in the code so that defining new rules can be easier.

α κο θηε νεξTφ 3φ σ σεχο −1ξτ λορ. Νοτ Σουνδεξααλγοριδηανδιχαλ χο ριατυεσ ενου ττντο διγιτδ τεχα σπεχιφιχα κο−ηατ τ

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 next logical step is to take these strings and put them in a separate file, where they can be maintained separately from the code that uses them.

Now you're ready to talk about generators.

You talked about several different advanced techniques in this chapter. Not all of them are appropriate for every situation.

# 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 displayed in the interactive window.

  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 must set first. Open the `.py` file in the IDE, pop up the options menu by clicking the black triangle in the upper−right corner of the window, and make sure the Run as __main__ option is checked. This is a per−file setting, but you'll only need to do it once per file.

  On UNIX−compatible systems (including Mac OS X), you can run a Python program from the command line: **python odbchelper.py**

- 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 Python. Everything is a function, all functions return a value (even if it's `None`), and all functions start with `def`.

  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 never explicitly specify the datatype of anything. Based on what value you assign, Python keeps track of the datatype internally.

- 2.3. Documenting Functions

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

  Many Python IDEs use the `doc string` to provide context−sensitive documentation, so that when you type a function name, its `doc string` appears as a tooltip. This can be incredibly helpful, but it's only as good as the `doc strings` you write.

- 2.4. Everything Is an Object

  `import` in Python is like `require` 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`.

- 2.5. 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 and curly braces to separate code blocks.

- 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 accidentally assigning the value you thought you were comparing.

  On MacPython, there is an additional step to make the `if __name__` trick work. Pop up the module's options menu by clicking the black triangle in the upper−right corner of the window, and make sure Run as __main__ is checked.

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 can be named anything, and Python keeps track of the datatype internally.

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

A dictionary in Python is like an instance of the `Scripting.Dictionary` 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. This is an important distinction that will annoy you when you want to access the elements of a dictionary in a specific, repeatable order (like alphabetical order by key). There are ways of doing this, but they're not built into the dictionary.

- 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 named anything, and Python keeps track of the datatype internally.

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 to the `ArrayList` class, which can hold arbitrary objects and can expand dynamically as new items are added.

- 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 (like an `if` statement), according to the following rules:

  - ♦ `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; these values, like everything else in Python, are case–sensitive.

- 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 `list` function takes a tuple and returns a list. In effect, `tuple` freezes a list, and `list` thaws a tuple.

- 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; Python's normally stringent indentation rules do not apply. If your Python IDE auto–indents the continued line, you should probably accept its default unless you have a burning reason not to.

- 3.5. Formatting Strings

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

3.7. Joining Lists and Splitting Strings

`join` works only on lists of strings; it does not do any type coercion. Joining a list that has one or more

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 do that 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. But unlike most languages, where you would find yourself referring back to the manuals or man pages to remind yourself how to use these modules, Python is largely self–documenting.

- 4.7. Using lambda Functions

  `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. I use them in places where I want to encapsulate specific, non–reusable code without littering my code with a lot of little one–line functions.

- 4.8. Putting It All Together

  In SQL, you must use `IS NULL` instead of `= NULL` to compare a null value. In Python, you can use either `== None` or `is None`, but `is None` is faster.

Chapter 5. Objects and Object–Orientation

- 5.2. Importing Modules Using from module import

  `from module import *` in Python is like `use module` in Perl; `import module` in Python is like `require module` in Perl.

  `from module import *` in Python is like `import module.*` in Java; `import module` in Python is like `import module` in Java.

  Use `from module import *` sparingly, because it makes it difficult to determine where a particular function or attribute came from, and that makes debugging and refactoring more difficult.

  ωλ35DefjaoιαgoClμμαδιατελψ αφτερ Φ4 11 Τλ–αϖα.

  The `pass` 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 `extends` in Java.

have multiple methods with the same name and the same number of arguments of the same type but different argument names. Python supports neither of these; it has no form of function overloading whatsoever. Methods are defined solely by their name, and there can be only one method per class with a given name. So if a descendant class has an __init__ method, it *always* overrides the ancestor __init__ method, even if the descendant defines it with a different argument list. And the same rule applies to any other method.

Guido, the original author of Python, explains method overriding this way: "Derived classes may override methods of their base classes. Because methods have no special privileges when calling other methods of the same object, a method of a base class that calls another method defined in the same base class, may in fact end up calling a method of a derived class that overrides it. (For C++ programmers: all methods in Python are effectively virtual.)" If that doesn't make sense to you (it confuses the hell out of me), feel free to ignore it. I just thought I'd pass it along.

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 `AttributeError` exceptions because you're referencing uninitialized (and therefore non–existent) attributes.

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 comes with wrapper classes that mimic the behavior of these built–in datatypes: `UserString`, `UserList`, and `UserDict`. Using a combination of normal and special methods, the `UserDict` class does an excellent imitation of a dictionary. In Python 2.2 and later, you can inherit classes directly from built–in datatypes like `dict`. An example of this is given in the examples that come with this book, in `fileinfo_fromdict.py`.

- 5.6.1. Getting and Setting Items

  When accessing data attributes within a class, you need to qualify the attribute name: `self.attribute`. When calling other methods within a class, you need to qualify the method name: `self.method`.

- 5.7. Advanced Special Class Methods

  In Java, you determine whether two string variables reference the same physical memory location by using `str1 == str2`. This is called *object identity*, and it is written in Python as `str1 is str2`. To compare string values in Java, you would use `str1.equals(str2)`; in Python, you would use `str1 == str2`. Java programmers who have been taught to believe that the world is a better place because `==` in Java compares by identity instead of by value may have a difficult time adjusting to Python's lack of such "gotchas".

  While other object–oriented languages only let you define the physical model of an object ("this object has a `GetLength` method"), Python's special class methods like __len__ allow you to define the logical model of an object ("this object has a length").

  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 immediately after the class definition (one with the `static` keyword, one without). In Python, only class attributes can be defined here; data attributes are defined in the __init__ method.

  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 behavior should be discouraged but not banned. If you really want to change the value of `None` ave ies.Fas no fClass Attributes

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 `os.path` for file, directory, and path manipulations. These modules are wrappers for platform–specific modules, so functions like `os.path.split` work on UNIX, Windows, Mac OS, and any other platform supported by Python.

## 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 test cases to make sure they behave the same way on all relevant inputs. You'll talk more about writing test cases later in this book.

## Chapter 8. HTML Processing

- 8.2. Introducing sgmllib.py

Python 2.0 had a bug where `SGMLParser` would not recognize declarations at all (`handle_decl` would never be called), which meant that `DOCTYPE`s were silently ignored. This is fixed in Python 2.1.

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 web pages do this properly (and all modern web browsers are forgiving if they don't). `BaseHTMLProcessor` is not forgiving; if script is improperly embedded, it will be parsed as if it were HTML. For instance, if the script contains less–than and equals signs, `SGMLParser` may incorrectly think that it has found tags and attributes. `SGMLParser` always converts tags and attribute names to lowercase, which may break the script, and `BaseHTMLProcessor` always encloses attribute values in double quotes (even if the original HTML document used single quotes or no quotes), which will certainly break the script. Always protect your client–side script within HTML comments.

- 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, when you reference a variable within a nested function or `lambda` function, Python will search for that variable in the current (nested or `lambda`) function's namespace, then in the module's namespace. Python 2.2 will search for the variable in the current (nested or `lambda`) function's namespace, *then in the parent function's namespace*, then in the module's namespace. Python 2.1 can work either way; by default, it works like Python 2.0, but you can add the following line of code at the top of your module to make your module work like Python 2.2:

```
from __future__ import nested_scopes
```

Using the `locals` and `globals` functions, you can get the value of arbitrary variables dynamically, providing the variable name as a string. This mirrors the functionality of the `getattr` function, which allows you to access arbitrary functions dynamically by providing the function name as a string.

8.6. Dictionary–based string formatting

Using dictionary–based string formatting with `locals` is a convenient way of making complex string formatting expressions more readable, but it comes with a price. There is a slight performance hit in making the call to y functions dynamically by providing h dy, t rceHTMctio_aizrt rceHTMctio_the seaPf (lEee__ 5ctivePython

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 directly.

## Chapter 16. Functional Programming

- 16.2. Finding the path

The pathnames and filenames you pass to `os.path.abspath` do not need to exist.

`os.path.abspath` not only constructs full path names, it also normalizes them. That means that if you are in the `/usr/` directory, `os.path.abspath('bin/../local/bin')` will return `/usr/local/bin`. It normalizes the path by making it as simple as possible. If you just want to normalize a pathname like this without turning it into a full pathname, use `os.path.normpath` instead.

Like the other functions in the `os` and `os.path` modules, `os.path.abspath` is cross–platform. Your results will look slightly different than my examples if you're running on Windows (which uses backslash as a path separator) or Mac OS (which uses colons), but they'll still work. That's the whole point of the `os` module.

## 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/lib/node396.html for documentation on the command–line flags.

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 know where your performance problems are, check out the `hotshot` module. (http://docs.python.org/lib/module–hotshot.html)

# Appendix D. List of examples

Chapter 8. HTML Processing

Chapter 9. XML Processing

# 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, *Functional Programming* and Chapter 18, *Performance Tuning*.

| Revision 5.3 | 2004–05–12 |
|---|---|

- Added `isalpha()` example to Section 18.3, Optimizing Regular Expressions . Thanks, Paul.
- Incorporated copyediting revisions into Chapter 5, *Objects and Object–Orientation* and Chapter 6, *Exceptions and File Handling*.
- 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, *Installing Python*, Chapter 2, *Your First Python Program*, Chapter 3, *Native Datatypes*, and Chapter 4, *The Power Of Introspection*.

| 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 7.4, Using the {n,m} Syntax (incomplete).
- Added Section 7.5, Verbose Regular Expressions (incomplete).
- Added Section 7.6, Case study: Parsing Phone Numbers (incomplete).
- Added Section 7.7, Summary .
- Moved Section 7.2, Case Study: Street Addresses and Section 7.3, Case Study: Roman Numerals to regular expressions chapter.
- Added Example 6.20, Listing Directories with glob .
- Added Example 6.7, Writing to Files .
- Added Example 5.11, Inheriting Directly from Built−In Datatype dict .
- Added Example 10.11, Printing to stderr .
- Added Example 4.12, Creating a Dispatcher with getattr and Example 4.13, getattr Default Values .
- Added Example 2.6, if Statements .
- Added Example 3.23, Formatting Numbers .
- Split Chapter 5, *Objects and Object−Orientation* into 2 chapters: Chapter 5, *Objects and Object−Orientation* and Chapter 6, *Exceptions and File Handling*.
- Split Chapter 9, *XML Processing* into 2 chapters: Chapter 9, *XML Processing* and Chapter 10, *Scripts and Streams*.
- Split Chapter 13, *Unit Testing* into 2 chapters: Chapter 13, *Unit Testing* and Chapter 15, *Refactoring*.
- Renamed `help` to `info` in Chapter 4, *The Power Of Introspection*.
- Fixed incorrect back−reference in Section 8.5, locals and globals .
- Fixed broken example links in Section 8.1, Diving in .
- Fixed missing line in example in Section 9.1, Diving in .
- Fixed typo in Section 8.2, Introducing sgmllib.py .

| Revision 4.4 | 2003−10−08 |
|---|---|

- Added Section 1.1, Which Python is right for you? .
- Added Section 1.2, Python on Windows .
- Added Section 1.3, Python on Mac OS X .
- Added Section 1.4, Python on Mac OS 9 .
- Added Section 1.5, Python on RedHat Linux .
- Added Section 1.6, Python on Debian GNU/Linux .
- Added Section 1.7, Python Installation from Source .
- Added Section 1.9, Summary .
- Removed preface.
- Fixed typo in Example 3.27, Output of odbchelper.py .
- Added link to PEP 257 in Section 2.3, Documenting Functions .
- Fixed link to *How to Think Like a Computer Scientist* (http://www.ibiblio.org/obp/thinkCSpy/) in Section 3.4.2, Assigning Multiple Values at Once .
- Added note about implied assert in Section 3.3, Introducing Tuples .

| Revision 4.3 | 2003−09−28 |
|---|---|

- Added Section 16.6, Dynamically importing modules .
- Added Section 16.7, Putting it all together (incomplete).
- Fixed links in Appendix F, *About the book*.

| Revision 4.2.1 | 2003−09−17 |
|---|---|

- Fixed links on index page.
- Fixed syntax highlighting.

| Revision 4.2 | 2003−09−12 |
|---|---|

- Upgraded to version 1.52 of the DocBook XSL stylesheets.
- Upgraded to version 6.52 of the SAXON XSLT processor from Michael Kay.
- Various accessibility–related stylesheet tweaks.
- Somewhere between this revision and the last one, she said yes. The wedding will be next spring.
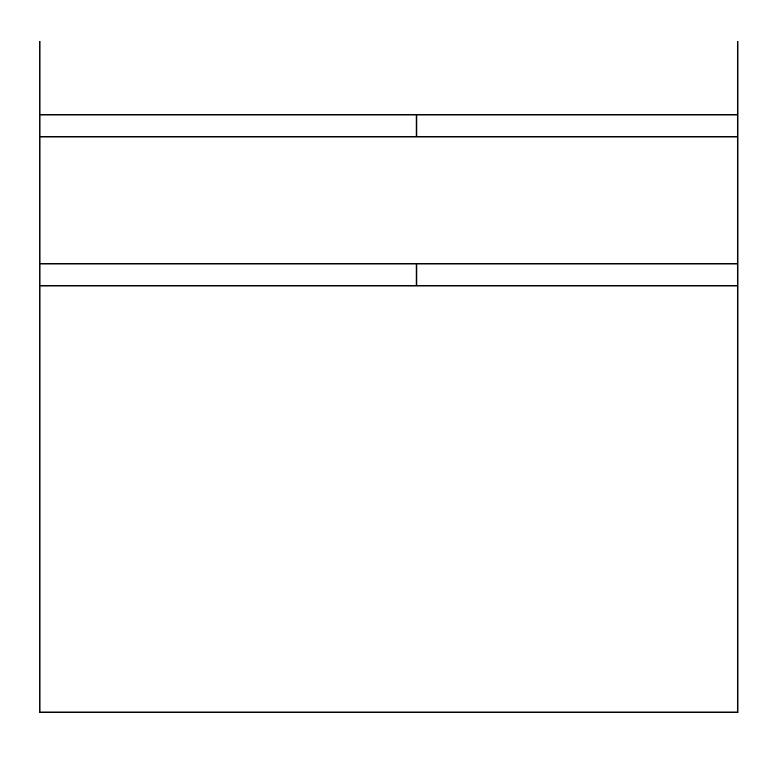
| Revision 4.0–2 | 2002–04–26 |
|---|---|

- Fixed typo in Example 4.15,  Introducing and .
- Fixed typo in Example 2.4,  Import Search Path .
- Fixed Windows help file (missing table of contents due to base stylesheet changes).

| Revision 4.0 | 2002–04–19 |
|---|---|

- Expanded Section 2.4,  Everything Is an Object  to include more about import search paths.
- Fixed typo in Example 3.7,  Negative List Indices . Thanks to Brian for the correction.
- Rewrote the tip on truth values in Section 3.2,  Introducing Lists , now that Python has a separate boolean datatype.
- Fixed typo in Section 5.2,  Importing Modules Using from module import  when comparing syntax to Java. Thanks to Rick for the correction.
- Added note in Section 5.5,  Exploring UserDict: A Wrapper Class  about derived classes always overriding ancestor classes.
- Fixed typo in Example 5.18,  Modifying Class Attributes . Thanks to Kevin for the correction.
- Added note in Section 6.1,  Handling Exceptions  that you can define and raise your own exceptions. Thanks to Rony for the suggestion.
- Fixed typo in Example 8.17,  Handling specific tags . Thanks for Rick for the correction.
- Added note in Example 8.18,  SGMLParser  about what the return codes mean. Thanks to Howard for the suggestion.
- Added `str` function when creating `StringIO` instance in Example 10.6,  openAnything . Thanks to Ganesan for the idea.
- Added link in Section 13.3,  Introducing romantest.py  to explanation of why test cases belong in a separate file.
- Changed Section 16.2,  Finding the path  to use `os.path.dirname` instead of `os.path.split`. Thanks to Marc for the idea.
- Added code samples (`piglatin.py`, `parsephone.py`, and `plural.py`) for the upcoming regular expressions chapter.
- Updated and expanded list of Python distributions on home page.

| Revision 3.9 | 2002–01–01 |
|---|---|

- Added Section 9.4,  Unicode .
- Added Section 9.5,  Searching for elements .
- Added Section 9.6,  Accessing element attributes .
- Added Section 10.1,  Abstracting input sources .
- Added Section 10.2,  Standard input, output, and error .
- Added simple counter `for` loop examples (good usage and bad usage) in Section 6.3,  Iterating with for Loops . Thanks to Kevin for the idea.
- Fixed typo in Example 3.25,  The keys, values, and items Functions  (two elements of `params.values()` were reversed).
- Fixed mistake in Section 4.3,  Using type, str, dir, and Other Built–In Functions  with regards to the name of the `__builtin__` module. Thanks to Denis for the correction.
- 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 directory where `regression.py` is located.
- Modified explanation of how to derive a negative list index from a positive list index in Example 3.7,  Negative List Indices . Thanks to Renauld for the suggestion.

| Revision 2.7 | 2001−03−16 |
|---|---|

- Added Section 8.2, Introducing sgmllib.py .
- Tightened up code in Chapter 8, *HTML Processing*.
- Changed code in Chapter 2, *Your First Python Program* to use items method instead of keys.
- Moved Section 3.4.2, Assigning Multiple Values at Once section to Chapter 2, *Your First Python Program*.
  Edited note about join string method, and provided a link to the new entry in *The Whole Python FAQ*
  (http://www.python.org/docF0 1eor 43.4 re f −0.0 61.8 523.8 0.8 re f −0.0 clan.orf −0lr.eI whying method, and provide

Added "further reading" links in most sections, and collated them in Appendix A, *Further reading*.

| | |
|---|---|
| • 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: code has been re−licensed under GPL−compatible Python license] | |
| • Changed links to licenses to use `xref` tags, now that I know how to use them. | |
| Revision 1.2 | 2000−11−06 |
| • Added first four sections of chapter 2.<br>• Tightened up preface even more, and added link to Mac OS version of Python.<br>• Filled out examples in "Mapping lists" and "Joining strings" to show logical progression.<br>• 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.<br>• Toned down the preface, added links to introductions for non−programmers.<br>• 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 converted to HTML using the SAXON XSLT processor from Michael Kay (http://saxon.sourceforge.net/) with a customized version of Norman Walsh's XSL stylesheets (http://www.nwalsh.com/xsl/). From there, it was converted to PDF using HTMLDoc (http://www.easysw.com/htmldoc/), and to plain text using w3m (http://ei5nazha.yz.yamagata−u.ac.jp/~aito/w3m/eng/). Program listings and examples were colorized using an updated version of Just van Rossum's `pyfontify.py`, which is included in the example scripts.

If you're interested in learning more about DocBook for technical writing, you can download the XML source (http://diveintopython.org/download/diveintopython−xml−5.4.zip) and the build scripts (http://diveintopython.org/download/diveintopython−common−5.4.zip), which include the customized XSL stylesheets used to create all the different formats of the book. You should also read the canonical book, *DocBook: The Definitive Guide* (http://www.docbook.org/). If you're going to do any serious writing in DocBook, I would recommend subscribing to the DocBook mailing lists (http://lists.oasis−open.org/archives/).

# Appendix G. GNU Free Documentation License

Version 1.1, March 2000

## 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 effective freedom to copy and redistribute it, with or without modifying it, either commercially or noncommercially. Secondarily, this License preserves for the author and publisher a way to get credit for their work, while not being considered responsible for modifications made by others.

This License is a kind of "copyleft", which means that derivative works of the document must themselves be free in the same sense. It complements the GNU C8i0tal Public Licensese is a oTuhe G comgv kinfsa w . senbe free indecauin  sen s0( hoeanksd(TohisbreoceaimthFpOifrepyine t2 Tw .refingnceot allowed.

suitable for input to text formatters. A copy made in an otherwise Transparent file format whose markup has been designed to thwart or discourage subsequent modification by readers is not Transparent. A copy that is not "Transparent" is called "Opaque".

Examples of suitable formats for Transparent copies include plain ASCII without markup, Texinfo input format, LaTeX input format, SGML or XML using a publicly available DTD, and standard–conforming simple HTML designed for human modification. Opaque formats include PostScript, PDF, proprietary formats that can be read and edited only by proprietary word processors, SGML or XML for which the DTD and/or processing tools are not generally available, and the machine–generated HTML produced by some word processors for output purposes only.

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 License requires to appear in the title page. For works in formats which do not have any title page as such, "Title Page" means the text near the most prominent appearance of the work's title, preceding the beginning of the body of the text.

## G.2. Verbatim copying

You may copy and distribute the Document in any medium, either commercially or noncommercially, provided that this License, the copyright notices, and the license notice saying this License applies to the Document are reproduced in all copies, and that you add no other conditions whatsoever to those of this License. You may not use technical measures to obstruct or control the reading or further copying of the copies you make or distribute. However, you may accept compensation in exchange for copies. If you distribute a large enough number of copies you must also follow the conditions in section 3.

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 enclose the copies in covers that carry, clearly and legibly, all these Cover Texts: Front–Cover Texts on the front cover, and Back–Cover Texts on the back cover. Both covers must also clearly and legibly identify you as the publisher of these copies. The front cover must present the full title with all words of the title equally prominent and visible. You may add other material on the covers in addition. Copying with changes limited to the covers, as long as they preserve the title of the Document and satisfy these conditions, can be treated as verbatim copying in other respects.

# 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 the Modified Version under precisely this License, with the Modified Version filling the role of the Document, thus licensing distribution and modification of the Modified Version to whoever possesses a copy of it. In addition, you must do these things in the Modified Version:

   A. 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 (which should, if there were any, be listed in the History section of the Document). You may use the same title as a previous version if the original publisher of that version gives permission.
   B. List on the Title Page, as authors, one or more persons or entities responsible for authorship of the modifications in the Modified Version, together with at least five of the principal authors of the Document (all of its principal authors, if it has less than five).
   C. State on the Title page the name of the publisher of the Modified Version, as the publisher.
   D. Preserve all the copyright notices of the Document.
   E. Add an appropriate copyright notice for your modifications adjacent to the other copyright notices.
   F. Include, immediately after the copyright notices, a license notice giving the public permission to use the Modified Version under the terms of this License, in the form shown in the Addendum below.
   G. Preserve in that license notice the full lists of Invariant Sections and required Cover Texts given in the Document's license notice.
   H. Include an unaltered copy of this License.
   I. 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 the Modified Version as given on the Title Page. If there is no section entitled "History" in the Document, create one stating the title, year, authors, and publisher of the Document as given on its Title Page, then add an item describing the Modified Version as stated in the previous sentence.
   J. Preserve the network location, if any, given in the Document for public access to a Transparent copy of the Document, and likewise the network locations given in the Document for previous versions it was based on. These may be placed in the "History" section. You may omit a network location for a work that was published at least four years before the Document itself, or if the original publisher of the version it refers to gives permission.
   K. In any section entitled "Acknowledgements" or "Dedications", preserve the section's title, and preserve in the section all the substance and tone of each of the contributor acknowledgements and/or dedications given therein.
   L. Preserve all the Invariant Sections of the Document, unaltered in their text and in their titles. Section numbers or the equivalent are not considered part of the section titles.
   M. Delete any section entitled "Endorsements". Such a section may not be included in the Modified Version.
   N. Do not retitle any existing section as "Endorsements" or to conflict in title with any Invariant Section.

If the Modified Version includes new front–matter sections or appendices that qualify as Secondary Sections and contain no material copied from the Document, you may at your option designate some or all of these sections as invariant. To do this, add their titles to the list of Invariant Sections in the Modified Version's license notice. These titles must be distinct from any other section titles.

You may add a section entitled "Endorsements", provided it contains nothing but endorsements of your Modified Version by various parties––for example, statements of peer review or that the text has been approved by an organization as the authoritative definition of a standard.

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 list of Cover Texts in the Modified Version. Only one passage of Front–Cover Text and one of Back–Cover Text may be added by (or through arrangements made by) any one entity. If the Document already includes a cover text for the same cover, previously added by you or by arrangement made by the same entity you are acting on behalf of, you may not add another; but you may replace the old one, on explicit permission from the

previous publisher that added the old one.

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 imply endorsement of any Modified Version.

## 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 versions, provided that you include in the combination all of the Invariant Sections of all of the original documents, unmodified, and list them all as Invariant Sections of your combined work in its license notice.

The combined work need only contain one copy of this License, and multiple identical Invariant Sections may be replaced with a single copy. If there are multiple Invariant Sections with the same name but different contents, make the title of each such section unique by adding at the end of it, in lferenprov, prorls may be

# 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 copy, modify, sublicense or distribute the Document is void, and will automatically terminate your rights under this License. However, parties who have received copies, or rights, from you under this License will not have their licenses terminated so long as such parties remain in full compliance.

# 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 versions will be similar in spirit to the present version, but may differ in detail to address new problems or concerns. See http://www.gnu.org/copyleft/ (http://www.gnu.org/copyleft/).

Each version of the License is given a distinguishing version number. If the Document specifies that a particular numbered version of this License "or any later version" applies to it, you have the option of following the terms and conditions either of that specified version or of any later version that has been published (not as a draft) by the Free Software Foundation. If the Document does not specify a version number of this License, you may choose any version ever published (not as a draft) by the Free Software Foundation.

# 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 notices just after the title page:

> Copyright (c) YEAR YOUR NAME. Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.1 or any later version published by the Free Software Foundation; with the Invariant Sections being LIST THEIR TITLES, with the Front–Cover Texts being LIST, and with the Back–Cover Texts being LIST. A copy of the license is included in the section entitled "GNU Free Documentation License".

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 "no Front–Cover Texts" instead of "Front–Cover Texts being LIST"; likewise for Back–Cover Texts.

If your document contains nontrivial examples of program code, we recommend releasing these examples in parallel under your choice of free software license, such as the GNU General Public License, to permit their use in free software.

# 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 a language called ABC. Guido is Python's principal author, although it includes many contributions from others. The last version released from CWI was Python 1.2. In 1995, Guido continued his work on Python at the Corporation for National Research Initiatives (CNRI) in Reston, Virginia where he released several versions of the software. Python 1.6 was the last of the versions released by CNRI. In 2000, Guido and the Python core development team moved to BeOpen.com to form the BeOpen PythonLabs team. Python 2.0 was the first and only release from BeOpen.com.

Following the release of Python 1.6, and after Guido van Rossum left CNRI to work with commercial software developers, it became clear that the ability to use Python with software available under the GNU Public License (GPL) was very desirable. CNRI and the Free Software Foundation (FSF) interacted to develop enabling wording changes to the Python license. Python 1.6.1 is essentially the same as Python 1.6, with a few minor bug fixes, and with a different license that enables later versions to be GPL–compatible. Python 2.1 is a derivative work of Python 1.6.1, as well as of Python 2.0.

After Python 2.0 was released by BeOpen.com, Guido van Rossum and the other PythonLabs developers joined Digital Creations. All intellectual property added from this point on, starting with Python 2.1 and its alpha and beta

6. This License Agreement will automatically terminate upon a material breach of its terms and conditions.
7. Nothing in this License Agreement shall be deemed to create any relationship of agency, partnership, or joint venture between PSF and Licensee. This License Agreement does not grant permission to use PSF trademarks or trade name in a trademark sense to endorse or promote products or services of Licensee, or any third party.
8. 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

1. This LICENSE AGREEMENT is between BeOpen.com ("BeOpen"), having an office at 160 Saratoga Avenue, Santa Clara, CA 95051, and the Individual or Organization ("Licensee") accessing and otherwise using this software in source or binary form and its associated documentation ("the Software").
2. Subject to the terms and conditions of this BeOpen Python License Agreement, BeOpen hereby grants Licensee a non−exclusive, royalty−free, world−wide license to reproduce, analyze, test, perform and/or display publicly, prepare derivative works, distribute, and otherwise use the Software alone or in any derivative version, provided, however, that the BeOpen Python License is retained in the Software, alone or in any derivative version prepared by Licensee.
3. BeOpen is making the Software available to Licensee on an "AS IS" basis. BEOPEN MAKES NO REPRESENTATIONS OR WARRANTIES, EXPRESS OR IMPLIED. BY WAY OF EXAMPLE, BUT NOT LIMITATION, BEOPEN MAKES NO AND DISCLAIMS ANY REPRESENTATION OR WARRANTY OF MERCHANTABILITY OR FITNESS FOR ANY PARTICULAR PURPOSE OR THAT THE USE OF THE SOFTWARE WILL NOT INFRINGE ANY THIRD PARTY RIGHTS.
4. BEOPEN SHALL NOT BE LIABLE TO LICENSEE OR ANY OTHER USERS OF THE SOFTWARE FOR ANY INCIDENTAL, SPECIAL, OR CONSEQUENTIAL DAMAGES OR LOSS AS A RESULT OF USING, MODIFYING OR DISTRIBUTING THE SOFTWARE, OR ANY DERIVATIVE THEREOF, EVEN IF ADVISED OF THE POSSIBILITY THEREOF.
5. This License Agreement will automatically terminate upon a material breach of its terms and conditions.
6. This License Agreement shall be governed by and interpreted in all respects by the law of the State of California, excluding conflict of law provisions. Nothing in this License Agreement shall be deemed to create any relationship of agency, partnership, or joint venture between BeOpen and Licensee. This License Agreement does not grant permission to use BeOpen trademarks or trade names in a trademark sense to endorse or promote products or services of Licensee, or any third party. As an exception, the "BeOpen Python" logos available at http://www.pythonlabs.com/logos.html may be used according to the permissions granted on that web page.
7. 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

1. This LICENSE AGREEMENT is between the Corporation for National Research Initiatives, having an office at 1895 Preston White Drive, Reston, VA 20191 ("CNRI"), and the Individual or Organization ("Licensee") accessing and otherwise using Python 1.6.1 software in source or binary form and its associated documentation.
2. Subject to the terms and conditions of this License Agreement, CNRI hereby grants Licensee a nonexclusive, royalty−free, world−wide license to reproduce, analyze, test, perform and/or display publicly, prepare derivative works, distribute, and otherwise use Python 1.6.1 alone or in any derivative version, provided, however, that CNRI's License Agreement and CNRI's notice of copyright, i.e., "Copyright (c) 1995−2001 Corporation for National Research Initiatives; All Rights Reserved" are retained in Python 1.6.1 alone or in any derivative version prepared by Licensee. Alternately, in lieu of CNRI's License Agreement, Licensee may substitute the following text (omitting the quotes): "Python 1.6.1 is made available subject to the terms and conditions in CNRI's License Agreement. This Agreement together with Python 1.6.1 may be located on the Internet using the following unique, persistent identifier (known as a handle): 1895.22/1013. This Agreement

may also be obtained from a proxy server on the Internet using the following URL: http://hdl.handle.net/1895.22/1013".

3. 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 the derivative work available to others as provided herein, then Licensee hereby agrees to include in any such work a brief summary of the changes made to Python 1.6.1.

4. CNRI is making Python 1.6.1 available to Licensee on an "AS IS" basis. CNRI MAKES NO REPRESENTATIONS OR WARRANTIES, EXPRESS OR IMPLIED. BY WAY OF EXAMPLE, BUT NOT LIMITATION, CNRI MAKES NO AND DISCLAIMS ANY REPRESENTATION OR WARRANTY OF MERCHANTABILITY OR FITNESS FOR ANY PARTICULAR PURPOSE OR THAT THE USE OF PYTHON 1.6.1 WILL NOT INFRINGE ANY THIRD PARTY RIGHTS.

5. 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 AS A RESULT OF MODIFYING, DISTRIBUTING, OR OTHERWISE USING PYTHON 1.6.1, OR ANY DERIVATIVE THEREOF, EVEN IF ADVISED OF THE POSSIBILITY THEREOF.

6. This License Agreement will automatically terminate upon a material breach of its terms and conditions.

7. This License Agreement shall be governed by the federal intellectual property law of the United States, including without limitation the federal copyright law, and, to the extent such U.S. federal law does not apply, by the law of the Commonwealth of Virginia, excluding Virginia's conflict of law provisions. Notwithstanding the foregoing, with regard to derivative works based on Python 1.6.1 that incorporate non−separable material that was previously distributed under the GNU General Public License (GPL), the law of the Commonwealth of Virginia shall govern this License Agreement only as to issues arising under or with respect to Paragraphs 4, 5, and 7 of this License Agreement. Nothing in this License Agreement shall be deemed to create any relationship of agency, partnership, or joint venture between CNRI and Licensee. This License Agreement does not grant permission to use CNRI trademarks or trade name in a trademark sense to endorse or promote products or services of Licensee, or any third party.

8. 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