

MASSIMO DI PIERRO

WEB AUTOMATION WITH PYTHON

EXPERTS4SOLUTIONS

Copyright 2008-2013 by Massimo Di Pierro. All rights reserved.

THE CONTENT OF THIS BOOK IS PROVIDED UNDER THE TERMS OF THE CREATIVE COMMONS PUBLIC LICENSE BY-NC-ND 3.0.

<http://creativecommons.org/licenses/by-nc-nd/3.0/legalcode>

THE WORK IS PROTECTED BY COPYRIGHT AND/OR OTHER APPLICABLE LAW. ANY USE OF THE WORK OTHER THAN AS AUTHORIZED UNDER THIS LICENSE OR COPYRIGHT LAW IS PROHIBITED.

BY EXERCISING ANY RIGHTS TO THE WORK PROVIDED HERE, YOU ACCEPT AND AGREE TO BE BOUND BY THE TERMS OF THIS LICENSE. TO THE EXTENT THIS LICENSE MAY BE CONSIDERED TO BE A CONTRACT, THE LICENSOR GRANTS YOU THE RIGHTS CONTAINED HERE IN CONSIDERATION OF YOUR ACCEPTANCE OF SUCH TERMS AND CONDITIONS.

Limit of Liability/Disclaimer of Warranty: While the publisher and author have used their best efforts in preparing this book, they make no representations or warranties with respect to the accuracy or completeness of the contents of this book and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives or written sales materials. The advice and strategies contained herein may not be suitable for your situation. You should consult with a professional where appropriate. Neither the publisher nor author shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

For more information about appropriate use of this material contact:

Massimo Di Pierro
School of Computing
DePaul University
243 S Wabash Ave
Chicago, IL 60604 (USA)
Email: massimo.dipierro@gmail.com

Library of Congress Cataloging-in-Publication Data:

ISBN: XXX-XXX-XXXX

Build Date: October 9, 2014

XXXXXXXXXXXX

Contents

1	Introduction	11
2	Linux	13
2.1	Access	13
2.2	Bash Shell	14
2.3	Piping stdin and stdout	16
2.4	Installing prepackaged apps	16
2.5	Emacs in a nutshell	17
2.6	Installing binary packages	18
2.7	Git Tutorial	18
2.7.1	Importing a new project	19
2.7.2	Viewing project history	21
2.7.3	Review of GIT commands	21
3	Python	23
3.0.4	Python vs Java and C++ syntax	25
3.0.5	help, dir	25
3.1	Python Types	26
3.1.1	int and long	27
3.1.2	float and decimal	27
3.1.3	str	30
3.1.4	list and array	31
3.1.5	tuple	33
3.1.6	dict	35
3.1.7	set	36

3.2	Python control flow statements	38
3.2.1	for...in	38
3.2.2	while	40
3.2.3	if...elif...else	41
3.2.4	try...except...else...finally	41
3.2.5	def...return	44
3.2.6	lambda	46
3.3	Classes	47
3.3.1	Special methods and operator overloading	49
3.4	File input/output	50
3.5	import modules	51
3.5.1	math and cmath	51
3.5.2	os	52
3.5.3	sys	53
3.5.4	datetime	53
3.5.5	time	54
3.5.6	urllib and json	54
3.6	Regular Expressions	55
4	OS funtions from Python	57
5	BeautifulSoup	61
5.1	Cached Downloads	63
5.2	A simple web crawler	63
6	Mechanize	67
7	Fabric	71
8	User Management	75
9	EC2 and OpenStack	77
9.1	EC2 Commad Line API Tools	78
	Bibliography	81

Preface

Some of the text in these notes has been collected from Wikipedia and other online sources.

1

Introduction

Your account:

```
1 > ssh user<yourid>@140.192.30.237
```


2

Linux

2.1 Access

In this course you will be using a Linux OS. There are many ways to use Linux. You can install it on your machine, with or without dual boot (dual boot is only recommended for experts), you can run it within a virtual machine (VMWare or VirtualBox for example), or you can access a remote Linux server.

Assuming you have a Linux account on a remote server you can connect to it using one of the following ways:

- From another Linux/Unix/Mac terminal:

```
1 > ssh <username>@<hostname>
```

for example:

```
1 > ssh thomas.anderson@the.matrix.net
```

- From Mac you can follow the same directions as Linux, but first you have to open Terminal. Search for “Terminal” in “Spotlight”.
- From Windows you need to download PuTTY from <http://the.earth.li/~sgtatham/putty/latest/x86/putty.exe> and use the PuTTY GUI to insert the username and hostname.
- If you have a PythonAnywhere.com account, login there, then click on

[Consoles] and [Start a new console][Bash]. PythonAnywhere provides a free Linux account in the cloud and they offer a Bash console shell you can access via the web browser.

2.2 Bash Shell

The “Bash shell” is a command line processor. It allows you to interactively type and execute commands and runs inside a text window. There are many different shells, a.k.a. command line processors and they have a different syntax. For example the Bash shell understand “bash” commands; the Python shell understand “Python” commands. Even Windows has a shell (the “Command Prompt”) which understand Windows commands. For now you will be using the Bash shell.

Example of commands:

What is my name?

```
1 > whoami
```

Make a file called “test.py” and write “print 1” in it:

```
1 > echo "print 1" > test.py
```

List all files in the current folder

```
1 > ls
```

(* means everything and it is referred to as a wildcard). List all files ending in .py

```
1 > ls *.py
```

List all files ending in .py with more details

```
1 > ls -l *.py
```

Run the Python program test.py

```
1 > python test.py
```

Delete all files ending in

```
1 > rm *~
```

List the name of the current folder

```
1 > pwd
```

Navigate to the parent of the current folder

```
1 > cd ..
```

Navigate to the top folder (also known as file system root)

```
1 > cd /
```

Navigate back to your home folder

```
1 > cd ~
```

Make a new subfolder

```
1 > mkdir work
```

Enter the subfolder

```
1 > cd work
```

Create an empty file called "test.txt"

```
1 > touch test.txt
```

Step out of a folder

```
1 > cd ..
```

Remove a folder and all its content

```
1 > rm -r work
```

Print a list of all programs I am running

```
1 > ps
```

Print a list of all programs running on this machine by or other user

```
1 > ps -aux
```

List all running programs and sort them by resource consumption

```
1 > top
```

(press "q" to quit) Start a program and send it in background

```
1 > python &
```

(the program will be running but will not IO to console) Kill a program currently running

```
1 > kill -9 <pid>
```

(here <pid> is the process id of the program as listed by top and or ps).

Practice with the above commands.

pwd, ls, cd, mkdir, echo, touch, rm, ps, top.

Use them to explore the file system.

Also notice that you can type an incomplete command and press [tab], the Bash shell will complete the command for you or give you a list of options. This is called autocomplete.

Some operations require superuser access and the system will refuse to perform the action otherwise. If you are an administrator of a system, you can perform an action with administrator by prepending `sudo` to any shell command.

For example - DO NOT DO THIS, EVER! - delete everything in the file system:

```
1 > cd /  
2 > sudo rm -r *
```

At this point you have the ability to create and browse files, folders, and running processes. A running program may be constituted of one or more processes (parallel programs may do more than one thing at once, thus more than one process).

2.3 Piping stdin and stdout

When using the `echo` command we used the following example:

```
1 > echo "print 1" > test.py
```

Here `>` is a redirection operator. It takes the output of the `echo` command which would normally print to console, and send it to the `test.py` file. Similarly, one can redirect standard input which is normally received from the keyboard, and read it instead from a file using the `<` operator.

2.4 Installing prepackaged apps

At this point we need to learn how to install new programs and run them. How you do this depends on the version of Linux. I will assume Ubuntu. I will also assume you have superuser access else you will not be able to install programs.

The easiest way to install prepackaged apps for Ubuntu is:

```
1 > sudo apt-get install <appname>
```

For example

```
1 > sudo apt-get install mc
```

Once installed we can start it with

```
1 > mc
```

mc is a useful extension to the Bash shell which allows us to type Bash commands and, at the same time, browse the file system in a visual way.

We can also re-install existing apps to make sure we have the latest version, for example:

```
1 > sudo apt-get install python
```

Other useful apps are the following

```
1 > sudo apt-get install zip      # for compressing files
2 > sudo apt-get install unzip   # for uncompressing
3 > sudo apt-get install tar     # for packing files
4 > sudo apt-get install emacs   # a file editor
5 > sudo apt-get install python  # the python interpreter
6 > sudo apt-get install g++     # the C++ compiler
7 > sudo apt-get install wget    # to download stuff from web
8 > sudo apt-get install make    # to install source packages
9 > sudo apt-get install build-essential
10 > sudo apt-get install python2.7 # to make sure we have it
11 > sudo apt-get install python-dev # useful development tools
12 > sudo apt-get install git     # version control system
```

2.5 Emacs in a nutshell

Emacs is a text editor. It understand the syntax of many programming languages and can help with development. It does not have a GUI but only a text based UI and you can perform various operations using command line shortcuts. Remember the most important of them all [ctrl]+g. if you get into trouble it brings you back to the edit interface.

You can open a file with:

```
1 > emacs test.py
```

[ctrl]+x+c	Save and Exit
[ctrl]+s	Search
[meta]+<	scroll to beginning
[meta]+>	scroll to end
[meta]+X	command prompt
[tab]	automatically re-indent the current line

2.6 Installing binary packages

Some time you may need to install binary packages. They come in the form of tar g-zipped files. They look like: `<appname>.tar.gz`

They installed using the following process:

```
1 > gunzip <appname>.tar.gz
2 > tar xvf <appname>.tar
3 > cd <appname-folder>
4 > ./configure
5 > make
6 > sudo make install
```

Here `<appname>` is the name of the application being installed. `<appname-folder>` is the folder created by `tar xvf ...`. The `tar` command extracts the content of the `.tar` file into a folder. The `./configure` step can take additional command line options and one should consult the documentation of the specific package for more details. `make` compiles the application. `make install` moves the compiled files in the proper place. If everything worked fine you can delete the `<appname-folder>`.

2.7 Git Tutorial

(from the Git man pages)

This tutorial explains how to import a new project into Git, make changes to it, and share changes with other developers.

If you are instead primarily interested in using Git to fetch a project, for

example, to test the latest version, you may prefer to start with the first two chapters of *The Git User's Manual*.

It is a good idea to introduce yourself to Git with your name and public email address before doing any operation. The easiest way to do so is:

```
1 > git config --global user.name "Your Name Comes Here"
2 > git config --global user.email you@yourdomain.example.com
```

2.7.1 Importing a new project

Assume you have a tar-ball `project.tar.gz` with your initial work. You can place it under Git revision control as follows.

```
1 > tar xzf project.tar.gz
2 > cd project
3 > git init
```

Git will reply with:

```
1 Initialized empty Git repository in .git/
```

You've now initialized the working directory—you may notice a new directory created, named `“.git”`.

Next, tell Git to take a snapshot of the contents of all files under the current directory (note the `.`), with `git add`:

```
1 > git add .
```

This snapshot is now stored in a temporary staging area which Git calls the `“index”`. You can permanently store the contents of the index in the repository with `git commit`:

```
1 > git commit
```

This will prompt you for a commit message. You've now stored the first version of your project in Git.

Making changes Modify some files, then add their updated contents to the index:

```
1 > git add file1 file2 file3
```

You are now ready to commit. You can see what is about to be committed using `git diff` with the `-cached` option:

```
1 > git diff --cached
```

(Without `-cached`, `git diff` will show you any changes that you've made but not yet added to the index.) You can also get a brief summary of the situation with `git status`:

```
1 > git status
2 # On branch master
3 # Changes to be committed:
4 #   (use `git reset HEAD <file>...' to unstage)
5 #
6 #modified:   file1
7 #modified:   file2
8 #modified:   file3
9 #
```

If you need to make any further adjustments, do so now, and then add any newly modified content to the index. Finally, commit your changes with:

```
1 > git commit
```

This will again prompt you for a message describing the change, and then record a new version of the project.

Alternatively, instead of running `git add` beforehand, you can use

```
1 > git commit -a
```

which will automatically notice any modified (but not new) files, add them to the index, and commit, all in one step.

A note on commit messages: Though not required, it's a good idea to begin the commit message with a single short (less than 50 character) line summarizing the change, followed by a blank line and then a more thorough description. The text up to the first blank line in a commit message is treated as the commit title, and that title is used throughout Git. For example, `git-format-patch(1)` turns a commit into email, and it uses the title on the Subject line and the rest of the commit in the body.

Git tracks content not files Many revision control systems provide an `add` command that tells the system to start tracking changes to a new file. Git's `add` command does something simpler and more powerful: `git add` is used both for new and newly modified files, and in both cases it takes a snapshot of the given files and stages that content in the index, ready for inclusion in the next commit.

2.7.2 Viewing project history

At any point you can view the history of your changes using

```
1 > git log
```

If you also want to see complete diffs at each step, use

```
1 > git log -p
```

Often the overview of the change is useful to get a feel of each step

```
1 > git log --stat --summary
```

Git is great for collaboration and allows creating and merging branches. When git projects are hosted on cloud services such as GitHub, you can send pull requests to other users which they can accept and reject. For more details consult the official documentation.

2.7.3 Review of GIT commands

3

Python

Python is a general-purpose high-level programming language. Its design philosophy emphasizes programmer productivity and code readability. It has a minimalist core syntax with very few basic commands and simple semantics. It also has a large and comprehensive standard library, including an Application Programming Interface (API) to many of the underlying operating system (OS) functions. Python provides built-in objects such as linked lists (`list`), tuples (`tuple`), hash tables (`dict`), arbitrarily long integers (`long`), complex numbers, and arbitrary precision decimal numbers.

Python supports multiple programming paradigms, including object-oriented (`class`), imperative (`def`), and functional (`lambda`) programming. Python has a dynamic type system and automatic memory management using reference counting (similar to Perl, Ruby, and Scheme).

Python was first released by Guido van Rossum in 1991[?]. The language has an open, community-based development model managed by the non-profit Python Software Foundation. There are many interpreters and compilers that implement the Python language, including one in Java (Jython), one built on .Net (IronPython) and one built in Python itself (PyPy). In this brief review, we refer to the reference C implementation created by Guido.

You can find many tutorials, the official documentation, and library references of the language on the official Python website. [1]

For additional Python references, we can recommend the books in ref. [?] and ref. [?].

You may skip this chapter if you are already familiar with the Python language.

3.0.4 Python vs Java and C++ syntax

	Java/C++	Python
assignment	$a = b;$	$a = b$
comparison	if ($a == b$)	if $a == b$:
loops	for($a = 0; a < n; a ++$)	for a in $range(0, n)$:
block	Braces {...}	indentation
function	$float f(float a) \{$	def $f(a)$:
function call	$f(a)$	$f(a)$
arrays/lists	$a[i]$	$a[i]$
member	$a.member$	$a.member$
nothing	$null / void*$	$None$

As in Java, variables which are primitive types (bool, int, float) are passed by copy but more complex types, unlike C++, are passed by reference.

3.0.5 help, dir

The Python language provides two commands to obtain documentation about objects defined in the current scope, whether the object is built-in or user-defined.

We can ask for help about an object, for example "1":

```

1 >>> help(1)
2 Help on int object:
3
4 class int(object)
5 |   int(x[, base]) -> integer
6 |
7 |   Convert a string or number to an integer, if possible. A floating point
8 |   argument will be truncated towards zero (this does not include a string
9 |   representation of a floating point number!) When converting a string, use
10 |   the optional base. It is an error to supply a base when converting a
11 |   non-string. If the argument is outside the integer range a long object
12 |   will be returned instead.
13 |
14 |   Methods defined here:
15 |
16 |   __abs__(...)
17 |       x.__abs__() <==> abs(x)

```

```
18 ...
```

and, since "1" is an integer, we get a description about the int class and all its methods. Here the output has been truncated because it is very long and detailed.

Similarly, we can obtain a list of methods of the object "1" with the command `dir`:

```
1 >>> dir(1)
2 ['__abs__', '__add__', '__and__', '__class__', '__cmp__', '__coerce__',
3  '__delattr__', '__div__', '__divmod__', '__doc__', '__float__',
4  '__floordiv__', '__getattr__', '__getnewargs__', '__hash__', '__hex__',
5  '__index__', '__init__', '__int__', '__invert__', '__long__', '__lshift__',
6  '__mod__', '__mul__', '__neg__', '__new__', '__nonzero__', '__oct__',
7  '__or__', '__pos__', '__pow__', '__radd__', '__rand__', '__rdiv__',
8  '__rdivmod__', '__reduce__', '__reduce_ex__', '__repr__', '__rfloordiv__',
9  '__rlshift__', '__rmod__', '__rmul__', '__ror__', '__rpow__', '__rrshift__',
10  '__rshift__', '__rsub__', '__rtruediv__', '__rxor__', '__setattr__',
11  '__str__', '__sub__', '__truediv__', '__xor__']
```

3.1 Python Types

Python is a dynamically typed language, meaning that variables do not have a type and therefore do not have to be declared. Variables may also change the type of value they hold through its life. Values, on the other hand, do have a type. You can query a variable for the type of value it contains:

```
1 >>> a = 3
2 >>> print(type(a))
3 <type 'int'>
4 >>> a = 3.14
5 >>> print(type(a))
6 <type 'float'>
7 >>> a = 'hello python'
8 >>> print(type(a))
9 <type 'str'>
```

Python also includes, natively, data structures such as lists and dictionaries.

3.1.1 int and long

There are two types representing integer numbers, they are `int` and `long`. The difference is `int` corresponds to the microprocessor's native bit length. Typically this is 32 bits and can hold signed integers in range $[-2^{31}, +2^{31}]$ while the `long` type can hold almost any arbitrary integer. It is important that Python automatically converts one into the other as necessary and you can mix and match the two types in computations. Here is an example:

```

1 >>> a = 1024
2 >>> type(a)
3 <type 'int'>
4 >>> b = a**128
5 >>> print(b)
6 20815864389328798163850480654728171077230524494533409610638224700807216119346720
7 59602447888346464836968484322790856201558276713249664692981627981321135464152584
8 82590187784406915463666993231671009459188410953796224233873542950969577339250027
9 68876520583464697770622321657076833170056511209332449663781837603694136444406281
10 042053396870977465916057756101739472373801429441421111406337458176
11 >>> print(type(b))
12 <type 'long'>

```

Computers represent 32 bit integer numbers by converting them to base 2. The conversion works in the following way:

```

1 def int2binary(n, nbits=32):
2     if n<0:
3         return [1 if bit==0 else 0 for bit in int2binary(-n-1,nbits)]
4     bits = [0]*nbits
5     for i in range(nbits):
6         n, bits[i] = divmod(n,2)
7     if n: raise OverflowError
8     return bits

```

The case $n < 0$ is called *two's complement* and is defined as the value obtained by subtracting the number from the largest power of two (2^{32} for 32 bits). Just by looking at the most significant bit one can determine the sign of the binary number (1 for negative and 0 for zero or positive).

3.1.2 float and decimal

There are two ways to represent decimal numbers in Python. Using the native double precision (64bits) representation, `float`, or using the `decimal` module.

Most numerical problems are dealt with simply using float:

```
1 >>> pi = 3.141592653589793
2 >>> two_pi = 2.0 * pi
```

The `cmath` module, described later, contains many math functions that use with the `float` type.

Floating point numbers are internally represented as follows:

$$x = \pm m2^e \quad (3.1)$$

where x is the number, m is called mantissa and it is zero or a number in range $[1,2)$, e is called exponent. The sign, m and e can be computed using the following algorithm which also writes their representation in binary:

```
1 def float2binary(x,nm=52,ne=11):
2     if x==0:
3         return 0, [0]*nm, [0]*ne
4     sign,mantissa, exponent = (1 if x<0 else 0),abs(x),0
5     while abs(mantissa)>=2:
6         mantissa,exponent = 0.5*mantissa,exponent+1
7     while 0<abs(mantissa)<1:
8         mantissa,exponent = 2.0*mantissa,exponent-1
9     return sign,int2binary(int(2**nm*(mantissa-1)),nm),int2binary(exponent,ne)
```

Because the mantissa is stored in a fixed number of bits (11 for a 64 bits floating point number) then exponents smaller than -1022 and larger than 1023 cannot be represented. An arithmetic operation that returns a number smaller than $2^{-1022} \simeq 10^{-308}$ cannot be represented and results in an Underflow error. An operation that returns a number larger than $2^{1023} \simeq 10^{308}$ also cannot be represented and results in an Overflow error.

Here is an example of Overflow:

```
1 >>> a = 10.0**200
2 >>> a*a
3 inf
```

And here is an example of Underflow:

```
1 >>> a = 10.0**-200
2 >>> a*a
3 0.0
```

Another problem with finite precision arithmetic is the loss of precision for

computation. Consider the case of the difference between two numbers with very different orders of magnitude. In order to compute the difference the CPU reduces them to the same exponent (the largest of the two) and then computes the difference in the two mantissas. If two numbers differ for a factor 2^k than the mantissa of the smallest number, in binary, needs to be shifted by k position thus resulting in a loss of information because the k least significant bits in the mantissa are ignored. If the difference between the two numbers is greater than a factor 2^{52} all bits in the mantissa of the smallest number are ignored and the smallest number becomes completely invisible.

Below is a practical example that produces a wrong result:

```

1 >>> a = 1.0
2 >>> b = 2.0**53
3 >>> a+b-b
4 0.0

```

Simple example of what occurs internally in a processor to add two floating point numbers together. The IEEE 754 standard states for 32 bit floating point numbers the exponent has a range of -126 to +127.

```

1 262 in IEEE 754: 0 10000111 000001100000000000000000 (+ e:8 m:1.0234375)
2 3 in IEEE 754: 0 10000000 100000000000000000000000 (+ e:1 m:1.5)
3 265 in IEEE 754: 0 10000111 000010010000000000000000

```

To add 262.0 to 3.0 the Exponents must be the same. The exponent of the lesser number is increased to the exponent of the greater number. In this case 3's exponent must be increased by 7. Increasing the exponent by 7 means the mantissa must be shifted 7 binary digits to the right.

```

1 0 10000111 000001100000000000000000
2 0 10000111 000000110000000000000000 (The implied '1' is also pushed 7 places to the
   right)
3 -----
4 0 10000111 000010010000000000000000 which is the IEEE 754 format for 265.0

```

In the case of two numbers where the exponent is greater than the number of digits in the mantissa the smaller number is shifted right off the end. The effect is a zero is added to the larger number.

In some cases only some of the bits of the smaller number's mantissa are lost an a partial addition occurs.


```

1 >>> print('number is ' + str(3))
2 number is 3
3 >>> print('number is %s' % (3))
4 number is 3
5 >>> print('number is %(number)s' % dict(number=3))
6 number is 3

```

The final notation is more explicit and less error prone, and is to be preferred.

Many Python objects, for example numbers, can be serialized into strings using `str` or `repr`. These two commands are very similar but produce slightly different output. For example:

```

1 >>> for i in [3, 'hello']:
2     ...     print(str(i), repr(i))
3 3 3
4 hello 'hello'

```

For user-defined classes, `str` and `repr` can be defined/redefined using the special operators `__str__` and `__repr__`. These are briefly described later in this chapter. For more information on the topic, refer to the official Python documentation [?].

Another important characteristic of a Python string is that it is an iterable object, similar to a list:

```

1 >>> for i in 'hello':
2     ...     print(i)
3 h
4 e
5 l
6 l
7 o

```

3.1.4 list **and** array

A list is a native mutable Python object made of a sequence of arbitrary objects. Arrays are not native to Python but unlike lists each element must be of the same type but are computationally faster.

The main methods of a Python lists and arrays are `append`, `insert`, and `delete`. Other useful methods include `count`, `index`, `reverse` and `sort`. Note: arrays do not have a `sort` method.

```

1 >>> b = [1, 2, 3]

```

```

2 >>> print(type(b))
3 <type 'list'>
4 >>> b.append(8)
5 >>> b.insert(2, 7) # insert 7 at index 2 (3rd element)
6 >>> del b[0]
7 >>> print(b)
8 [2, 7, 3, 8]
9 >>> print(len(b))
10 4
11 >>> b.append(3)
12 >>> b.reverse()
13 print(b, " 3 appears ", b.count(3), " times. The number 7 appears at index ", b.
    index(7))
14 [3, 8, 3, 7, 2] 3 appears 2 times. The number 7 appears at index 3

```

Lists can be sliced:

```

1 >>> a = [2, 7, 3, 8]
2 >>> print(a[:3])
3 [2, 7, 3]
4 >>> print(a[1:])
5 [7, 3, 8]
6 >>> print(a[-2:])
7 [3, 8]

```

and concatenated/joined:

```

1 >>> a = [2, 7, 3, 8]
2 >>> a = [2, 3]
3 >>> b = [5, 6]
4 >>> print(a + b)
5 [2, 3, 5, 6]

```

A list is iterable; you can loop over it:

```

1 >>> a = [1, 2, 3]
2 >>> for i in a:
3     ...     print(i)
4 1
5 2
6 3

```

There is a very common situation for which a *list comprehension* can be used. Consider the following code:

```

1 >>> a = [1,2,3,4,5]
2 >>> b = []
3 >>> for x in a:
4     ...     if x % 2 == 0:
5     ...         b.append(x * 3)
6 >>> print(b)

```



```
7 [6, 12]
```

This code clearly processes a list of items, selects and modifies a subset of the input list, and creates a new result list. This code can be entirely replaced with the following list comprehension:

```
1 >>> a = [1,2,3,4,5]
2 >>> b = [x * 3 for x in a if x % 2 == 0]
3 >>> print(b)
4 [6, 12]
```

Python has a module called `array`. It provides an efficient array implementation. Unlike lists, array elements must all be of the same type and the type must be either a char, short, int, long, float or double. A type of char, short, int or long may be either signed or unsigned.

```
1 >>> from array import array
2 >>> a = array('d', [1,2,3,4,5])
3 array('d', [1.0, 2.0, 3.0, 4.0, 5.0])
```

An array object can be used in the same way as a list but its elements must all be of the same type, specified by the first argument of the constructor (“d” for double, “l” for signed long, “f” for float, and “c” for character). For a complete list of available options, refer to the official Python documentation.

Using “array” over “list” can be faster but, more importantly, the “array” storage is more compact for large arrays.

3.1.5 tuple

A tuple is similar to a list, but its size and elements are immutable. If a tuple element is an object, the object itself is mutable but the reference to the object is fixed. A tuple is defined by elements separated by a comma and optionally delimited by round brackets:

```
1 >>> a = 1, 2, 3
2 >>> a = (1, 2, 3)
```

The round brackets are required for a tuple of zero elements such as

```
1 >>> a = () # this is an empty tuple
```

A trailing comma is required for a one element tuple but not for two or more elements.

34 WEB AUTOMATION WITH PYTHON

```
1 >>> a = (1) # not a tuple
2 >>> a = (1,) # this is a tuple of one element
3 >>> b = (1,2) # this is a tuple of two elements
```

Since lists are mutable, this works:

```
1 >>> a = [1, 2, 3]
2 >>> a[1] = 5
3 >>> print(a)
4 [1, 5, 3]
```

the element assignment does not work for a tuple:

```
1 >>> a = (1, 2, 3)
2 >>> print(a[1])
3 2
4 >>> a[1] = 5
5 Traceback (most recent call last):
6   File "<stdin>", line 1, in <module>
7 TypeError: 'tuple' object does not support item assignment
```

A tuple, like a list, is an iterable object. Notice that a tuple consisting of a single element must include a trailing comma, as shown below:

```
1 >>> a = (1)
2 >>> print(type(a))
3 <type 'int'>
4 >>> a = (1,)
5 >>> print(type(a))
6 <type 'tuple'>
```

Tuples are very useful for efficient packing of objects because of their immutability. The brackets are often optional. You may easily get each element of a tuple by assigning multiple variables to a tuple at one time:

```
1 >>> a = (2, 3, 'hello')
2 >>> (x, y, z) = a
3 >>> print(x)
4 2
5 >>> print(z)
6 hello
7 >>> a = 'alpha', 35, 'sigma' # notice the rounded brackets are optional
8 >>> p, r, q = a
9 print(r)
10 35
```

3.1.6 dict

A Python dictionary is a hash table that maps a key object to a value object. For example:

```

1 >>> a = {'k': 'v', 'k2': 3}
2 >>> print(a['k'])
3 v
4 >>> print(a['k2'])
5 3
6 >>> 'k' in a
7 True
8 >>> 'v' in a
9 False

```

You will notice that the format to define a dictionary is the same as JavaScript Object Notation [JSON]. Dictionaries may be nested:

```

1 >>> a = {'x': 3, 'y': 54, 'z': {'a': 1, 'b': 2}}
2 >>> print(a['z'])
3 {'a': 1, 'b': 2}
4 >>> print(a['z']['a'])
5 1

```

Keys can be of any hashable type (int, string, or any object whose class implements the `__hash__` method). Values can be of any type. Different keys and values in the same dictionary do not have to be of the same type. If the keys are alphanumeric characters, a dictionary can also be declared with the alternative syntax:

```

1 >>> a = dict(k='v', h2=3)
2 >>> print(a['k'])
3 v
4 >>> print(a)
5 {'h2': 3, 'k': 'v'}

```

Useful methods are `has_key`, `keys`, `values`, `items` and `update`:

```

1 >>> a = dict(k='v', k2=3)
2 >>> print(a.keys())
3 ['k2', 'k']
4 >>> print(a.values())
5 [3, 'v']
6 >>> a.update({'n1': 'new item'}) # adding a new item
7 >>> a.update(dict(n2='newer item')) # alternate method to add a new item
8 >>> a['n3'] = 'newest item' # another method to add a new item
9 >>> print(a.items())

```

```
10 [(('k2', 3), ('k', 'v'), ('n3', 'newest item'), ('n2', 'newer item'), ('n1', 'new
    item'))]
```

The `items` method produces a list of tuples, each containing a key and its associated value.

Dictionary elements and list elements can be deleted with the command `del`:

```
1 >>> a = [1, 2, 3]
2 >>> del a[1]
3 >>> print(a)
4 [1, 3]
5 >>> a = dict(k='v', h2=3)
6 >>> del a['h2']
7 >>> print(a)
8 {'k': 'v'}
```

Internally, Python uses the hash operator to convert objects into integers, and uses that integer to determine where to store the value. Using a key that is not hashable will cause an `unhashable type error`

```
1 >>> hash("hello world")
2 -1500746465
3 >>> k = [1,2,3]
4 >>> a = {k:'4'}
5 Traceback (most recent call last):
6   File "<stdin>", line 1, in <module>
7 TypeError: unhashable type: 'list'
```

3.1.7 set

A set is something between a list and dictionary. It represents a non-ordered list of unique elements. Elements in a set cannot be repeated. Internally it is implemented as a hash-table, similar to a set of keys in a dictionary. A set is created using the set constructor. Its argument can be a list, a tuple, or an iterator:

```
1 >>> s = set([1,2,3,4,5,5,5]) # notice duplicate elements are removed
2 >>> print(s)
3 set([1,2,3,4,5])
4 >>> s = set((1,2,3,4,5))
5 >>> print(s)
6 set([1,2,3,4,5])
7 >>> s = set(i for i in range(1,6))
8 >>> print(s)
9 set([1, 2, 3, 4, 5])
```

Since sets are non-ordered list, appending to the end is not applicable. Instead of append, add elements to a set using the add method:

```
1 >>> s = set()
2 >>> s.add(2)
3 >>> s.add(3)
4 >>> s.add(2)
5 >>> print(s)
6 set([2, 3])
```

Notice that the same element can not be added twice (2 in the example). There is no exception/error thrown when trying to add the same element more than once.

Since sets are non-ordered, the order you add items is not necessarily the order they will be returned.

```
1 >>> s = set([6, 'b', 'beta', -3.4, 'a', 3, 5.3])
2 >>> print(s)
3 set(['a', 3, 6, 5.3, 'beta', 'b', -3.4])
```

The set object supports normal set operations like union, intersection, and difference:

```
1 >>> a = set([1,2,3])
2 >>> b = set([2,3,4])
3 >>> c = set([2,3])
4 >>> print(a.union(b))
5 set([1, 2, 3, 4])
6 >>> print(a.intersection(b))
7 set([2, 3])
8 >>> print(a.difference(b))
9 set([1])
10 >>> if len(c) == len(a.intersection(c)):
11 ...     print("c is a subset of a")
12 ... else:
13 ...     print("c is not a subset of a")
14 ...
15 c is a subset of a
```

to check for membership:

```
1 >>> 2 in a
2 True
```

3.2 Python control flow statements

Python uses indentation to delimit blocks of code. A block starts with a line ending with colon and continues for all lines that have a similar or higher indentation as the next line. For example:

```

1 >>> i = 0
2 >>> while i < 3:
3 ...     print(i)
4 ...     i = i + 1
5 0
6 1
7 2

```

It is common to use four spaces for each level of indentation. It is a good policy not to mix tabs with spaces, which can result in (invisible) confusion.

3.2.1 for...in

In Python, you can loop over iterable objects:

```

1 >>> a = [0, 1, 'hello', 'python']
2 >>> for i in a:
3 ...     print(i)
4 0
5 1
6 hello
7 python

```

In the example above you will notice that the loop index, 'i', takes on the values of each element in the list [0, 1, 'hello', 'python'] sequentially. Python range keyword creates a list of integers automatically, that may be used in a 'for' loop without manually creating a long list of numbers.

```

1 >>> a = range(0,5)
2 >>> print(a)
3 [0, 1, 2, 3, 4]
4 >>> for i in a:
5 ...     print(i)
6 0
7 1
8 2
9 3
10 4

```

The parameters for `range(a,b,c)` are first parameter is the starting value of the list. The second parameter is next value if the list contained one more element. The third parameter is the increment value.

`range` can also be called with one parameter. It is matched to "b" above with the first parameter defaulting to 0 and the third to 1.

```

1 >>> print(range(5))
2 [0, 1, 2, 3, 4]
3 >>> print(range(53,57))
4 [53,54,55,56]
5 >>> print(range(102,200,10))
6 [102, 112, 122, 132, 142, 152, 162, 172, 182, 192]
7 >>> print(range(0,-10,-1))
8 [0, -1, -2, -3, -4, -5, -6, -7, -8, -9]
```

`range` is very convenient for creating a list of numbers, however as the list grows in length, the memory required to store the list also grows. A more efficient option is to use the keyword `xrange`. It generates an iterable range instead of generating the entire list of elements.

```

1 >>> for i in xrange(0, 4):
2     print(i)
3 0
4 1
5 2
6 3
```

This is equivalent to the C/C++/C#/Java syntax:

```

1 for(int i=0; i<4; i=i+1) { print(i); }
```

Another useful command is `enumerate`, which counts while looping and returns a tuple consisting of (index, value):

```

1 >>> a = [0, 1, 'hello', 'python']
2 >>> for (i, j) in enumerate(a): # the ( ) around i, j are optional
3     print(i, j)
4 0 0
5 1 1
6 2 hello
7 3 python
```

There is also a keyword `range(a, b, c)` that returns a list of integers starting with the value `a`, incrementing by `c`, and ending with the last value smaller than `b`, `a` defaults to 0 and `c` defaults to 1. `xrange` is similar but does not actually generate the list, only an iterator over the list; thus it is better for

looping.

You can jump out of a loop using `break`

```
1 >>> for i in [1, 2, 3]:
2     ...     print(i)
3     ...     break
4 1
```

You can jump to the next loop iteration without executing the entire code block with `continue`

```
1 >>> for i in [1, 2, 3]:
2     ...     print(i)
3     ...     continue
4     ...     print('test')
5 1
6 2
7 3
```

Python also supports list comprehensions and you can build lists using using the following syntax:

```
1 >>> a = [i*i for i in [0, 1, 2, 3]:
2 >>> print(a)
3 [0, 1, 4, 9]
```

Sometimes you may need a counter to “count” the elements of a list while looping:

```
1 >>> a = [e*(i+1) for (i,e) in ['a','b','c','d']]
2 >>> print(a)
3 ['a', 'bb', 'ccc', 'dddd']
```

3.2.2 while

Comparison operators in Python follow the C/C++/Java operators of `==`, `!=`, ... etc. However Python also accepts the `<>` operator as not equal to and is equivalent to `!=`. Logical operators are `and`, `or` and `not`.

The `while` loop in Python works much as it does in many other programming languages, by looping an indefinite number of times and testing a condition before each iteration. If the condition is `False`, the loop ends.

```
1 >>> i = 0
2 >>> while i < 10:
3     ...     i = i + 1
```



```

4 >>> print(i)
5 10

```

The for loop was introduced earlier in this chapter.

There is no loop...until or do...while construct in Python.

3.2.3 if...elif...else

The use of conditionals in Python is intuitive:

```

1 >>> for i in range(3):
2 ...     if i == 0:
3 ...         print('zero')
4 ...     elif i == 1:
5 ...         print('one')
6 ...     else:
7 ...         print('other')
8 zero
9 one
10 other

```

elif means "else if". Both elif and else clauses are optional. There can be more than one elif but only one else statement. Complex conditions can be created using the not, and and or logical operators.

```

1 >>> for i in range(3):
2 ...     if i == 0 or (i == 1 and i + 1 == 2):
3 ...         print('0 or 1')

```

3.2.4 try...except...else...finally

Python can throw - pardon, raise - Exceptions:

```

1 >>> try:
2 ...     a = 1 / 0
3 ... except Exception, e:
4 ...     print('oops: %s' % e)
5 ... else:
6 ...     print('no problem here')
7 ... finally:
8 ...     print('done')
9 oops: integer division or modulo by zero
10 done

```

42 WEB AUTOMATION WITH PYTHON

If an exception is raised, it is caught by the `except` clause and the `else` clause is not not executed. The `finally` clause is always executed.

There can be multiple `except` clauses for different possible exceptions:

```
1 >>> try:
2 ...     raise SyntaxError
3 ... except ValueError:
4 ...     print('value error')
5 ... except SyntaxError:
6 ...     print('syntax error')
7 syntax error
```

The `finally` clause is guaranteed to be executed while the `except` and `else` are not. In the example below the function returns within a `try` block. This is bad practice, but it shows that the `finally` will execute regardless of the reason the `try` block is exited.

```
1 >>> def f(x):
2 ...     try:
3 ...         r = x*x
4 ...         return r # bad practice
5 ...     except:
6 ...         print("exception occurred %s" % e)
7 ...     else:
8 ...         print("nothing else to do")
9 ...     finally:
10 ...        print("Finally we get here")
11 ...
12 >>> y = f(3)
13 Finally we get here
14 >>> print "result is ", y
15 result is 9
```

For every `try` you must have either a `except` or `finally` while the `else` is optional

Here is a list of built-in Python exceptions

```
1 BaseException
2 +-- SystemExit
3 +-- KeyboardInterrupt
4 +-- Exception
5     +-- GeneratorExit
6     +-- StopIteration
7     +-- StandardError
8         | +-- ArithmeticError
9         | | +-- FloatingPointError
```

```

10 | | +-- OverflowError
11 | | +-- ZeroDivisionError
12 | | +-- AssertionError
13 | | +-- AttributeError
14 | | +-- EnvironmentError
15 | | | +-- IOError
16 | | | +-- OSError
17 | | | | +-- WindowsError (Windows)
18 | | | | +-- VMSError (VMS)
19 | | +-- EOFError
20 | | +-- ImportError
21 | | +-- LookupError
22 | | | +-- IndexError
23 | | | +-- KeyError
24 | | +-- MemoryError
25 | | +-- NameError
26 | | | +-- UnboundLocalError
27 | | +-- ReferenceError
28 | | +-- RuntimeError
29 | | | +-- NotImplementedError
30 | | +-- SyntaxError
31 | | | +-- IndentationError
32 | | | | +-- TabError
33 | | +-- SystemError
34 | | +-- TypeError
35 | | +-- ValueError
36 | | | +-- UnicodeError
37 | | | | +-- UnicodeDecodeError
38 | | | | +-- UnicodeEncodeError
39 | | | | +-- UnicodeTranslateError
40 | +-- Warning
41 | | +-- DeprecationWarning
42 | | +-- PendingDeprecationWarning
43 | | +-- RuntimeWarning
44 | | +-- SyntaxWarning
45 | | +-- UserWarning
46 | | +-- FutureWarning
47 | | +-- ImportWarning
48 | | +-- UnicodeWarning

```

For a detailed description of each of them, refer to the official Python documentation.

Any object can be raised as an exception, but it is good practice to raise objects that extend one of the built-in exception classes.

44 WEB AUTOMATION WITH PYTHON

3.2.5 def...return

Functions are declared using `def`. Here is a typical Python function:

```
1 >>> def f(a, b):
2 ...     return a + b
3 >>> print(f(4, 2))
4 6
```

There is no need (or way) to specify the type of an argument(s) or the return value(s). In this example, a function `f` is defined that can take two arguments.

Functions are the first code syntax feature described in this chapter to introduce the concept of *scope*, or *namespace*. In the above example, the identifiers `a` and `b` are undefined outside of the scope of function `f`:

```
1 >>> def f(a):
2 ...     return a + 1
3 >>> print(f(1))
4 2
5 >>> print(a)
6 Traceback (most recent call last):
7   File "<pyshell#22>", line 1, in <module>
8     print(a)
9 NameError: name 'a' is not defined
```

Identifiers defined outside of function scope are accessible within the function; observe how the identifier `a` is handled in the following code:

```
1 >>> a = 1
2 >>> def f(b):
3 ...     return a + b
4 >>> print(f(1))
5 2
6 >>> a = 2
7 >>> print(f(1) # new value of a is used)
8 3
9 >>> a = 1 # reset a
10 >>> def g(b):
11 ...     a = 2 # creates a new local a
12 ...     return a + b
13 >>> print(g(2))
14 4
15 >>> print(a # global a is unchanged)
16 1
```

If `a` is modified, subsequent function calls will use the new value of the global `a` because the function definition binds the storage location of the identifier

`a`, not the value of `a` itself at the time of function declaration; however, if `a` is assigned-to inside function `g`, the global `a` is unaffected because the new local `a` hides the global value. The external-scope reference can be used in the creation of *closures*:

```

1 >>> def f(x):
2 ...     def g(y):
3 ...         return x * y
4 ...     return g
5 >>> doubler = f(2) # doubler is a new function
6 >>> tripler = f(3) # tripler is a new function
7 >>> quadrupler = f(4) # quadrupler is a new function
8 >>> print(doubler(5))
9 10
10 >>> print(tripler(5))
11 15
12 >>> print(quadrupler(5))
13 20

```

Function `f` creates new functions; and note that the scope of the name `g` is entirely internal to `f`. Closures are extremely powerful.

Function arguments can have default values and can return multiple results as a tuple: (Notice the parentheses are optional and are omitted in the example.)

```

1 >>> def f(a, b=2):
2 ...     return a + b, a - b
3 >>> x, y = f(5)
4 >>> print(x)
5 7
6 >>> print(y)
7 3

```

Function arguments can be passed explicitly by name, therefore the order of arguments specified in the caller can be different than the order of arguments with which the function was defined:

```

1 >>> def f(a, b=2):
2 ...     return a + b, a - b
3 >>> x, y = f(b=5, a=2)
4 >>> print(x)
5 7
6 >>> print(y)
7 -3

```

46 WEB AUTOMATION WITH PYTHON

Functions can also take a runtime-variable number of arguments. Parameters that start with `*` and `**` must be the last two parameters. If the `**` parameter is used it must be last in the list. Extra values passed in will be placed in the `*identifier` parameter while named values will be placed into the `ft**identifier`. Notice that when passing values into the function the unnamed values must be before any and all named values.

```
1 >>> def f(a, b, *extra, **extraNamed):
2 ...     print "a = ", a
3 ...     print "b = ", b
4 ...     print "extra = ", extra
5 ...     print "extranamed = ", extraNamed
6 >>> f(1, 2, 5, 6, x=3, y=2, z=6)
7 a = 1
8 b = 2
9 extra = (5, 6)
10 extranamed = {'y': 2, 'x': 3, 'z': 6}
```

Here the first two paramters (1 and 2) are matched with the paramters a and b while the tuple 5, 6 is placed into extra and the remaining items (which are in a dictionary format) are placed into extraNamed

In the opposite case, a list or tuple can be passed to a function that requires individual positional arguments by unpacking them:

```
1 >>> def f(a, b):
2 ...     return a + b
3 >>> c = (1, 2)
4 >>> print(f(*c))
5 3
```

and a dictionary can be unpacked to deliver keyword arguments:

```
1 >>> def f(a, b):
2 ...     return a + b
3 >>> c = {'a':1, 'b':2}
4 >>> print(f(**c))
5 3
```

3.2.6 lambda

lambda provides a way to define a short unnamed function:

```
1 >>> a = lambda b: b + 2
2 >>> print(a(3))
3 5
```

The expression "`lambda [a]:[b]`" literally reads as "a function with arguments [a] that returns [b]". The `lambda` expression is itself unnamed, but the function acquires a name by being assigned to identifier `a`. The scoping rules for `def` apply to `lambda` equally, and in fact the code above, with respect to `a`, is identical to the function declaration using `def`:

```
1 >>> def a(b):
2     ...     return b + 2
3 >>> print(a(3))
4 5
```

The only benefit of `lambda` is brevity; however, brevity can be very convenient in certain situations. Consider a function called `map` that applies a function to all items in a list, creating a new list:

```
1 >>> a = [1, 7, 2, 5, 4, 8]
2 >>> map(lambda x: x + 2, a)
3 [3, 9, 4, 7, 6, 10]
```

This code would have doubled in size had `def` been used instead of `lambda`. The main drawback of `lambda` is that (in the Python implementation) the syntax allows only for a single expression; however, for longer functions, `def` can be used and the extra cost of providing a function name decreases as the length of the function grows. Just like `def`, `lambda` can be used to *curry* functions: new functions can be created by wrapping existing functions such that the new function carries a different set of arguments:

```
1 >>> def f(a, b): return a + b
2 >>> g = lambda a: f(a, 3)
3 >>> g(2)
4 5
```

Python functions, created with either `def` or `lambda` allow re-factoring of existing functions in terms of a different set of arguments.

3.3 Classes

Because Python is dynamically typed, Python classes and objects may seem odd. In fact, member variables (attributes) do not need to be specifically defined when declaring a class and different instances of the same class can have different attributes. Attributes are generally associated with the

instance, not the class (except when declared as "class attributes", which is the same as "static member variables" in C++/Java).

Here is an example:

```

1 >>> class MyClass(object): pass
2 >>> myinstance = MyClass()
3 >>> myinstance.myvariable = 3
4 >>> print(myinstance.myvariable)
5 3

```

Notice that `pass` is a do-nothing command. In this case it is used to define a class `MyClass` that contains nothing. `MyClass()` calls the constructor of the class (in this case the default constructor) and returns an object, an instance of the class. The `(object)` in the class definition indicates that our class extends the built-in object class. This is not required, but it is good practice.

Here is a more involved class with multiple methods:

```

1 >>> class Complex(object):
2 ...     z = 2
3 ...     def __init__(self, real=0.0, imag=0.0):
4 ...         self.real, self.imag = real, imag
5 ...     def magnitude(self):
6 ...         return (self.real**2 + self.imag**2)**0.5
7 ...     def __add__(self, other):
8 ...         return Complex(self.real+other.real, self.imag+other.imag)
9 >>> a = Complex(1,3)
10 >>> b = Complex(2,1)
11 >>> c = a + b
12 >>> print(c.magnitude())
13 5

```

Functions declared inside the class are methods. Some methods have special reserved names. For example, `__init__` is the constructor. In the example we created a class to store the `real` and the `imag` part of a complex number. The constructor takes these two variables and stores them into `self` (not a keyword but a variable that plays the same role as `this` in Java and `(*this)` in C++. (This syntax is necessary to avoid ambiguity when declaring nested classes, such as a class that is local to a method inside another class, something the Python allows but Java and C++ do not).

The `self` variable is defined by the first argument of each method. They all must have it but they can use another variable name. Even if we use another

name, the first argument of a method always refers to the object calling the method. It plays the same role as the `this` keyword in the Java and the C++ languages.

`__add__` is also a special method (all special methods start and end in double underscore) and it overloads the `+` operator between `self` and `other`. In the example, `a+b` is equivalent to a call to `a.__add__(b)` and the `__add__` method receives `self=a` and `other=b`.

All variables are local variables of the method except variables declared outside methods which are called *class variables*, equivalent to C++ *static member variables* that hold the same value for all instances of the class.

3.3.1 Special methods and operator overloading

Class attributes, methods, and operators starting with a double underscore are usually intended to be private (*for example* to be used internally but not exposed outside the class) although this is a convention that is not enforced by the interpreter.

Some of them are reserved keywords and have a special meaning.

For example:

- `__len__`
- `__getitem__`
- `__setitem__`

They can be used, for example, to create a container object that acts like a list:

```

1 >>> class MyList(object):
2 >>>     def __init__(self, *a): self.a = list(a)
3 >>>     def __len__(self): return len(self.a)
4 >>>     def __getitem__(self, key): return self.a[key]
5 >>>     def __setitem__(self, key, value): self.a[key] = value
6 >>> b = MyList(3, 4, 5)
7 >>> print(b[1])
8 4
9 >>> b.a[1] = 7
10 >>> print(b.a)
11 [3, 7, 5]
```

Other special operators include `__getattr__` and `__setattr__`, which define the get and set methods (getters and setters) for the class, and `__add__`, `__sub__`, `__mul__`, `__div__` which overload arithmetic operators. For the use of these operators we refer the reader to the chapter on linear algebra where they will be used to implement algebra for matrices.

3.4 File input/output

In Python you can open and write in a file with:

```
1 >>> file = open('myfile.txt', 'w')
2 >>> file.write('hello world')
3 >>> file.close()
```

Similarly, you can read back from the file with:

```
1 >>> file = open('myfile.txt', 'r')
2 >>> print(file.read())
3 hello world
```

Alternatively, you can read in binary mode with "rb", write in binary mode with "wb", and open the file in append mode "a", using standard C notation.

The read command takes an optional argument, which is the number of bytes. You can also jump to any location in a file using seek.

You can read back from the file with read

```
1 >>> print(file.seek(6))
2 >>> print(file.read())
3 world
```

and you can close the file with:

```
1 >>> file.close()
```

In the standard distribution of Python, which is known as CPython, variables are reference-counted, including those holding file handles, so CPython knows that when the reference count of an open file handle decreases to zero, the file may be closed and the variable disposed. However, in other implementations of Python such as PyPy, garbage collection is used instead of reference counting, and this means that it is possible that there may accumulate too many open file handles at one time, resulting in an error before the gc has a chance to close and dispose of them all. Therefore it is

best to explicitly close file handles when they are no longer needed.

3.5 import modules

The real power of Python is in its library modules. They provide a large and consistent set of Application Programming Interfaces (APIs) to many system libraries (often in a way independent of the operating system).

For example, if you need to use a random number generator, you can do:

```
1 >>> import random
2 >>> print(random.randint(0, 9))
3 5
```

This prints a random integer in the range of (0,9], 5 in the example. The function `randint` is defined in the module `random`. It is also possible to import an object from a module into the current namespace:

```
1 >>> from random import randint
2 >>> print(randint(0, 9))
```

or import all objects from a module into the current namespace:

```
1 >>> from random import *
2 >>> print(randint(0, 9))
```

or import everything in a newly defined namespace:

```
1 >>> import random as myrand
2 >>> print(myrand.randint(0, 9))
```

In the rest of this book, we will mainly use objects defined in modules `math`, `cmath`, `os`, `sys`, `datetime`, `time` and `cPickle`. We will also use the `random` module but we will describe it in a later chapter.

In the following subsections we consider those modules that are most useful.

3.5.1 math and cmath

Here is a sampling of some of the methods available in the `math` and `cmath` packages.

- `math.isinf(x)` returns true if the floating point number `x` is positive or negative infinity
- `math.isnan(x)` returns true if the floating point number `x` is NaN. See Python documentation or IEEE 754 standards for more information.
- `math.exp(x)` returns e^x
- `math.log(x[, base])` returns the logarithm of `x` to the optional base. If `base` is not supplied `e` is assumed.
- `math.cos(x)`, `math.sin(x)`, `math.tan(x)` returns the cos, sin, tan of the value of `x`. `x` is in radians.
- `math.pi`, `math.e` the constants for `pi` and `e` to available precision

3.5.2 os

This module provides an interface to the operating system API. For example:

```
1 >>> import os
2 >>> os.chdir('.')
3 >>> os.unlink('filename_to_be_deleted')
```

Some of the `os` functions, such as `chdir`, are not thread-safe, *for example* they should not be used in a multi-threaded environment.

`os.path.join` is very useful; it allows the concatenation of paths in an OS-independent way:

```
1 >>> import os
2 >>> a = os.path.join('path', 'sub_path')
3 >>> print(a)
4 path/sub_path
```

System environment variables can be accessed via:

```
1 >>> print(os.environ)
```

which is a read-only dictionary.

3.5.3 sys

The `sys` module contains many variables and functions, but the used the most is `sys.path`. It contains a list of paths where Python searches for modules. When we try to import a module, Python searches the folders listed in `sys.path`. If you install additional modules in some location and want Python to find them, you need to append the path to that location to `sys.path`.

```
1 >>> import sys
2 >>> sys.path.append('path/to/my/modules')
```

3.5.4 datetime

The use of the `datetime` module is best illustrated by some examples:

```
1 >>> import datetime
2 >>> print(datetime.datetime.today())
3 2008-07-04 14:03:90
4 >>> print(datetime.date.today())
5 2008-07-04
```

Occasionally you may need to time-stamp data based on the UTC time as opposed to local time. In this case you can use the following function:

```
1 >>> import datetime
2 >>> print(datetime.datetime.utcnow())
3 2008-07-04 14:03:90
```

The `datetime` module contains various classes: `date`, `datetime`, `time` and `timedelta`. The difference between two date or two datetime or two time objects is a `timedelta`:

```
1 >>> a = datetime.datetime(2008, 1, 1, 20, 30)
2 >>> b = datetime.datetime(2008, 1, 2, 20, 30)
3 >>> c = b - a
4 >>> print(c.days)
5 1
```

We can also parse dates and datetimes from strings for example:

```
1 >>> s = '2011-12-31'
2 >>> a = datetime.datetime.strptime(s, '%Y-%m-%d') #modified
3 >>> print(s.year, s.day, s.month)
4 2011 31 12 #modified
```

54 WEB AUTOMATION WITH PYTHON

Notice that “%Y” matches the 4-digits year, “%m” matches the month as a number (1-12), “%d” matches the day (1-31), “%H” matches the hour, “%M” matches the minute, and “%S” matches the second. Check the Python documentation for more options.

3.5.5 time

The time module differs from date and datetime because it represents time as seconds from the epoch (beginning of 1970).

```
1 >>> import time
2 >>> t = time.time()
3 1215138737.571
```

Refer to the Python documentation for conversion functions between time in seconds and time as a datetime.

3.5.6 urllib and json

The urllib is a module to download data or a web page from a URL.

```
1 >>> page = urllib.urlopen('http://www.google.com/')
2 >>> html = page.read()
```

3.6 Regular Expressions

The most efficient way of searching for text/patterns in text is by using regular expressions. For example:

```

1 import urllib
2 import re
3 html = urllib.urlopen('http://www.depaul.com')
4 regex = re.compile('[\w\.\.]+@[\w\.\.]+')
5 emails = regex.findall(html)
6 print emails

```

`re.compile(...)` compiles the description of a pattern into an object, `regex`. The `findall` method of the object finds all expressions matching the compiled pattern. Here are useful rules:

- `^` means “beginning of text”.
- `$` means “and of text of text”.
- `.` means any character.
- `\w` means any alphanumeric character or underscore.
- `\W` means any character not alphanumeric character and not underscore.
- `\s` means any form of whitespace (space, tab, etc.).
- `\S` means any char but not a whitespace.
- `[...]` can be used to build an OR sequence of characters.
- `(...)` can be used to build a AND sequence of consecutive matching characters.
- `(...)*` means any repetition of the sequence `(...)`.
- `(...)+` means any non-zero repetition of the sequence `(...)`.

A slash can be used to escape special characters. In a string `\` must be escaped as `\\`. An actual “`\`” must therefore be escaped twice `\\\\`.

Here are some examples of regular expressions (they are all approximations and oversimplifications):

- Match all email in a document

```
1 re.compile("[\w\.\.]+@[ \w\.\. ]").findall(html)
```

- Match all tags in HTML

```
1 re.compile("\<(\w+)").findall(html)
```

- Match all images sources in HTML

```
1 re.compile("""\<img [^>]*src=[\''](.*)"[\''"]""").findall(html.lower())
```

- Match all links in HTML

```
1 re.compile("""\<a [^>]*href=[\''](.*)"[\''"]""").findall(html.lower())
```


4

OS functions from Python

Get current working directory with `os.getcwd()`

```
1 print os.getcwd()
```

Get the status of a file with `os.stat()`

```
1 print "Getting the status of: ", os.stat('/usr/bin/python')
```

Execute a shell command with `os.system()`

```
1 os.system('ls -l')
```

Return the current process id with `os.getpid()`

```
1 print os.getpid()
```

Change the owner and group id of path to the numeric uid and gid with `os.chown()`

```
1 os.chown(path, uid, gid)
```

Processes in the system run queue averaged over the last 1, 5, and 15 minutes

```
1 print os.getloadavg()
```

Check if a path exists with `os.path.exists()`

```
1 if os.path.exists("file.txt"):
```

Create a new directory named 'new_directory' if it doesn't exist already"

```
1 os.path.exists("new_directory") or os.mkdir("new_directory")
```

Check if the path is a directory or a file with `os.path.isdir()` and `os.path.isfile()`

```
1 path = "/tmp"
```

```

2 if os.path.isdir(path): print "That's a directory"
3 if os.path.isfile(path): print "That's a file"

```

Create a directory with `os.mkdir()`

```
1 print os.mkdir('new_directory', 0666)
```

Recursive create directories with `os.makedirs()`

```
1 os.makedirs('dir_a/dir_b/dir_c')
```

Remove a directory with `os.rmdir()`

```
1 print os.rmdir('directory')
```

Recursively remove empty directories with `os.removedirs()`

```
1 os.removedirs('dir_a/dir_b/dir_c')
```

Rename a file with `os.rename()`

```
1 print os.rename('/path/to/old/file', '/path/to/new/file')
```

Rename a file with `shutil.move()`

```
1 print shutil.move('/path/to/old/file', '/path/to/new/file')
```

Rename a file with `shutil.copy()`

```
1 print shutil.copy('/path/to/old/file', '/path/to/new/file')
```

Get the users home directory

```
1 print os.path.expanduser('~')
```

Check read permissions with `os.access()`

```

1 path = '/tmp/file.txt'
2 print os.access(path, os.R_OK)

```

Get the users environment with `os.environmen()`

```

1 home = os.environ['HOME']
2 print home

```

Move focus to a different directory with `os.chdir()`

```
1 print os.chdir('/tmp')
```

Print out all directories, sub-directories and files with `os.walk()`

```

1 for root, dirs, files in os.walk("/tmp"):
2     print root
3     print dirs
4     print files

```

Get the last time a directory was accessed with `os.path.getatime()`

```
1 os.path.getatime('/tmp')
```

Get the last time a directory was modified with `os.path.getmtime()`

```
1 os.path.getmtime('/tmp')
```

Get the user ID with `os.getuid()`

```
1 if os.getuid() != 0: print "you are not root"
```

Get the group ID with `os.getgid()`

```
1 print os.getgid()
```

Return the name of the user logged in with `os.getlogin()`

```
1 print os.getlogin()
```

Returns a list of all files in a directory with `os.listdir()`

```
1 for filename in os.listdir("/tmp"):
2     print "This is inside /tmp", filename
```

Get the size of a file with `os.path.getsize()`

```
1 os.path.getsize("/tmp/file.txt")
```

Using Python in your daily work is a good way to automate system administration tasks, when you feel that your shell scripts are too limited.

5

BeautifulSoup

Regular expressions are useful and indispensable, yet building the proper expressions to match generic tags, and tag attributes in HTML can be difficult. A better solution is to use a library designed specifically for this. One such library is BeautifulSoup. BeautifulSoup does not completely eliminate the need for regular expressions (in fact it uses regular expressions and you may have to build some to pass to some BeautifulSoup methods).

You can install BeautifulSoup from the Bash shell with

```
1 > sudo easy_install BeautifulSoup
```

From within Python you only need to import one class BeautifulSoup which we shall rename Soup for brevity:

```
1 >>> from BeautifulSoup import BeautifulSoup as Soup
```

Given any help text, for example:

```
1 >>> html = "<html><p>Para 1<p>Para 2<blockquote>Quote 1<blockquote>Quote 2"
```

You can use the Soup object to parse it and construct an object representation of the Document Object Model:

```
1 >>> soup = Soup(html)
```

this object can be serialized back into HTML (although it would not produce the exact same HTML, but an equivalent one):

```
1 >>> print soup.prettify()  
2 <html>
```

62 WEB AUTOMATION WITH PYTHON

```
3 <p>
4   Para 1
5 </p>
6 <p>
7   Para 2
8   <blockquote>
9     Quote 1
10    <blockquote>
11      Quote 2
12    </blockquote>
13  </blockquote>
14 </p>
15 </html>
```

The soup object (an instance of BeautifulSoup or BeautifulSoup) is a deeply-nested, well-connected data structure that corresponds to the structure of an XML or HTML document. The parser object contains two other types of objects: Tag objects, which correspond to tags like the <TITLE> tag and the tags; and NavigableString objects, which correspond to strings like “Page title” and “This is paragraph”.

The soup object can be used to query the DOM and manipulate it. For example, you can ask for the first <p/> tag which has an “align” attribute equal to “center”:

```
1 >>> p = soup.find('p', align='center')
```

The output p object is also a BeautifulSoup object. Unless you copy it, any change to this object will also affect the soup object that contains it.

Given then p object you can retrieve its attributes, for example its id, its class, its contents, or its serialized content:

```
1 >>> print p['id']
2 >>> print p['class']
3 >>> print p.contents
4 >>> print p.string
```

Similarly you can search all the tags in the soup object which match a given criteria, for example app <p/> tags:

```
1 >>> soup.findAll('p')
```

or all tags with an id that matches a regular expression:

```
1 >>> soup.findAll(id=re.compile("^X"))
```

An alternative syntax for getting all the <p/> tags is:

```
1 >>> tags = soup.p
```

The `'.'` notation can be chained to obtain, for example the title inside the head of the document:

```
1 >>> print soup.head.title
```

You can also replace its contents with:

```
1 >>> soup.head.title.replaceWith('New Title')
```

Each BeautifulSoup object is a reference to a portion of the soup document. Any change will affect the original document.

Notice that any HTML document is a tree, therefore you can traverse the tree by looking for the parent of a given node, its siblings (next and previous), as well as the next and previous tags in the DOM structure:

```
1 >>> print soup.head.parent.name
2 >>> print soup.p.nextSibling
3 >>> print soup.p.previousSibling
4 >>> print soup.p.next
5 >>> print soup.p.previous
```

5.1 Cached Downloads

```
1 import urllib2
2 import shelve
3
4 def cached_download(url):
5     d = shelve.open('cached.pages')
6     if url in d:
7         data = d[url]
8     else:
9         try:
10            data = d[url] = urllib2.urlopen(url, timeout=5).read()
11        except:
12            data = None
13    return data
```

5.2 A simple web crawler

```
1 # author: Massimo Di Pierro
2 # license: BSD
3 from BeautifulSoup import BeautifulSoup
4 import collections
```

64 WEB AUTOMATION WITH PYTHON

```
5 import cPickle
6 import urllib2
7 import shelve
8 import socket
9 import re
10
11 TIMEOUT = 10 #seconds
12 EXTENSIONS = ('asp','aspx','php','html','htm')
13 regex_remove = re.compile(
14     '|<script.*?</script>|<style.*?</style>|<!--.*?-->',re.DOTALL)
15 regex_words = re.compile('[a-zA-Z\-\_][a-zA-Z\-\_]+')
16 regex_url = re.compile('(P<scheme>\w+)\|:\/(P<host>.*?)(P<path_info>/.*?')
17
18 def cached_download(url):
19     d = shelve.open('cached.pages')
20     try:
21         data = d[url]
22     except KeyError:
23         try:
24             data = d[url] = urllib2.urlopen(url,timeout=TIMEOUT).read()
25         except (socket.gaierror, socket.error, AttributeError, IOError):
26             return None
27     return data
28
29 def parse_words(html):
30     body = html.find('body')
31     if not body: return set()
32     texts = body.findAll(text=True)
33     texts = map(lambda e:e.string.strip().lower(),texts)
34     texts = filter(lambda e:e,texts)
35     sets = [set(regex_words.findall(t)) for t in texts]
36     words = reduce(lambda a,b:a|b,sets,set())
37     return words
38
39 def normalize_url(new, scheme, host, path_info):
40     new = new.encode('utf8')
41     extension = new.split('?')[0].split('#')[0].rsplit('.')[-1].lower()
42     if (new.startswith('mailto:') or new.startswith('#') or
43         (not '/' in extension and not extension in EXTENSIONS) or
44         ('://' in new and not new.split('/:')[0] in ('http','https'))):
45         return None
46     elif '://' in new:
47         return new
48     elif new.startswith('//'):
49         return scheme+' '+new
50     elif new.startswith('/'):
51         return '%S://%S%S' % (scheme, host, new)
52     else:
53         relative = path_info.rsplit('/',1)[0]
```



```

54     while new.startswith('..'):
55         relative, new = relative.rsplit('/',1)[0], new[3:]
56     return '%s://%s%s/%s' % (scheme, host, relative, new)
57
58 def find_links(html, url = 'http://cdm.depaul.edu/', pattern=re.compile('.*')):
59     match = regex_url.match(url)
60     (scheme, host, path_info) = match.groups()
61     start = '/'.join(url.split('/')[1:3])
62     links = html.findAll('a')
63     links = map(lambda link:link.get('href'),links)
64     links = filter(lambda e:e,links)
65     links = map(lambda link:normalize_url(link,scheme,host,path_info),links)
66     follow_links = set(filter(lambda e:e and pattern.match(e),links))
67     ignore_links = set(filter(lambda e:e and e not in follow_links,links))
68     return follow_links, ignore_links
69
70 def crawl(urls = ['http://web2py.com/'], pattern = None,
71              max_links=200, filename='maps.pickle', mode='ignore-check'):
72     if isinstance(urls,str): urls = [urls]
73     maps = collections.defaultdict(set)
74     discovered = set(urls)
75     ignored = set()
76     queue = [url for url in urls]
77     n = 0
78     if pattern is None:
79         pattern = '^+'.join(re.escape(url) for url in urls)
80     if isinstance(pattern,str):
81         pattern = re.compile(pattern)
82     while queue and (max_links is None or n<max_links):
83         url = queue.pop()
84         print '%s: RETRIEVING %s' % (n, url)
85         data = cached_download(url)
86         if data is None:
87             print '%s: INVALID %s' % (n, url)
88             continue
89         html = BeautifulSoup(data)
90         follow_links, ignore_links = find_links(html,url,pattern)
91         words = parse_words(html)
92         for word in words:
93             maps[word].add(url)
94         if mode == 'ignore-log':
95             for url in ignore_links - ignored:
96                 print 'IGNORING %s' % url
97                 ignored |= ignore_links
98         elif mode == 'ignore-check':
99             for url in ignore_links - ignored:
100                 if cached_download(url) is None:
101                     print '%s: IGNORING BROKEN %s' % (n, url)
102                 else:

```

66 WEB AUTOMATION WITH PYTHON

```
103         print '%s: IGNORING %s' % (n, url)
104         ignored |= ignore_links
105     for url in follow_links:
106         if url and not url in discovered:
107             queue.append(url)
108             discovered.add(url)
109     print '%s: PROGRESS #urls:%s #words:%s' % (n, len(discovered), len(maps))
110     n += 1
111     cPickle.dump(maps, open(filename, 'w'))
112
113 def search(words, filename='maps.pickle'):
114     words = filter(lambda w: len(w) > 2, words.lower().split())
115     maps = cPickle.load(open(filename))
116     return reduce(lambda a, b: a&b, [maps[word] for word in words])
117
118 if __name__ == '__main__':
119     print crawl('http://depaul.edu/', pattern='^http://[\\w\\.]*depaul\\.edu(\\.*)',
120               max_links=None)
121     print search('Python programming')
```

6

Mechanize

You can install Mechanize from the Bash shell with

```
1 > sudo easy_install mechanize
```

Mechanize simulates a web browser and its interaction with the web. You can start it with:

```
1 >>> import mechanize
2 >>> br = mechanize.Browser()
```

It is important to configure your browser so that it can fool third party web servers. Web servers set restriction and who can access web pages and in particular they do not like to be visited by robots. You have to pretend your browser is a “normal” one, for example Firefox:

```
1 >>> br.set_handle_robots( False )
2 >>> br.addheaders = [('User-agent', 'Firefox')]
```

Once you have a browser `br` you can give instructions to it, for example, “visit the Google web page”:

```
1 >>> response = br.open("http://www.google.com/")
```

The response has attributes for example:

```
1 >>> print response.geturl()
2 http://www.google.com/
```

or the returned status code:

```
1 >>> print response.code
2 200
```

68 WEB AUTOMATION WITH PYTHON

or the response headers

```
1 >>> print response.info().headers
2 ['Date: Wed, 11 Sep 2013 18:09:10 GMT\r\n', 'Expires: -1\r\n', 'Cache-Control:
   private, max-age=0\r\n', 'Content-Type: text/html; charset=ISO-8859-1\r\n', '
   Set-Cookie: PREF=ID=a53b86f876c89fc9:U=548b6d330fff7040:FF=0:TM=1378922843:LM
   =1378922950:S=xPNYhTPc_HwsBqd8; expires=Fri, 11-Sep-2015 18:09:10 GMT; path=/;
   domain=.google.com\r\n', 'Server: gws\r\n', 'X-XSS-Protection: 1; mode=block\r\n',
   'X-Frame-Options: SAMEORIGIN\r\n', 'Alternate-Protocol: 80:quic\r\n', '
   Connection: close\r\n']
```

Or a specific header:

```
1 >>> response.info().getheader('date')
2 'Wed, 11 Sep 2013 18:09:10 GMT'
```

Or the page content (HTML or binary depending on the data):

```
1 >>> print response.read()
```

What is most important is the ability of Mechanize to understand the page. You can ask it for example to follow a link and get another page. For example, follow a link that contains text matching a regular expressions:

```
1 >>> response = br.follow_link(text_regex=".*Google.*")
```

You can ask if the action was successful and the browser received HTML:

```
1 >>> assert br.viewing_html()
```

Because of Mechanize ability to understand HTML, you can ask it to locate all links in a page:

```
1 >>> for f in br.links(): print f
2 ...
3 Link(base_url='http://www.google.com/', url='http://www.google.com/intl/en/options/
   ', text='More &#9660;', tag='a', attrs=[('class', 'gb3'), ('href', 'http://www
   .google.com/intl/en/options/'), ('onclick', 'this.blur();gbar.tg(event);return
   !1'), ('aria-haspopup', 'true')])
4 ...
```

And/or all forms in the page:

```
1 >>> for f in br.forms(): print f
2 <f GET http://www.google.ca/search application/x-www-form-urlencoded
3 <HiddenControl(ie=ISO-8859-1) (readonly)>
4 <HiddenControl(hl=en) (readonly)>
5 <HiddenControl(source=hp) (readonly)>
6 <TextControl(q=)>
7 <SubmitControl(btnG=Google Search) (readonly)>
8 <SubmitControl(btnI=I'm Feeling Lucky) (readonly)>
9 <HiddenControl(gbv=1) (readonly)>>
```

This says: There is only one form on the page/ Hidden field id's are ie (page encoding), hl (language code), hp (? don't know), and gbv (also don't know). The only not-hidden field id is q, which is a text input, which is the search text.

And you can ask it to select a form, fill it, and submit it:

```
1 >>> formnames = [f.name for f in br.forms()]
2 >>> br.select_form(formnames[0])
```

OR

```
1 >>> br.select_form(nr=0)
2 >>> br.form[ 'q' ] = 'foo'
3 >>> br.submit()
4 >>> response = br.response()
```


7

Fabric

Fabric is a remote management tool written in Python.

You can install Fabric from the Bash shell with

```
1 > sudo easy_install fabric
```

To use Fabric you create a fabfile (just a Python program) which defines the operations you want to perform for example:

```
1 from fabric.api import local
2
3 def dir():
4     local('ls -l')
```

than you run it with:

```
1 fab -f /path/to/fabfile.py dir
```

This calls the uptime function inside the "fabfile.py". From now on we will assume the "fabfile.py" in the user root folder and will omit the -f command line option.

The local(...) function executes the shell command passed as argument. In this case the ls -l command. The command in the example is executed on the local machine but the point of Fabric is to execute the command remotely and simultaneously on one or more machines. Here is a better example

```
1 from fabric.api import env, run
2
3 env.hosts = [ '192.168.1.100', '192.168.1.101', '192.168.1.102' ]
```

72 WEB AUTOMATION WITH PYTHON

```
4 env.user = "you"
5 env.password = "mysecret"
6 env.parallel = True
7
8 def dir():
9     run('ls -l')
```

Now a call to

```
1 fab dir
```

will produce a directory listing from all the machines at the hosts defined in `env.hosts`. One can also specify a host using the `-H` command line option:

```
1 fab -H 192.168.1.100 dir
```

You can use `-h` to find out more about command line options:

```
1 Usage: fab [options] <command>[:arg1,arg2=val2,host=foo,hosts='h1;h2',...] ...
2
3 Options:
4 -h, --help            show this help message and exit
5 -d NAME, --display=NAME
6                       print detailed info about command NAME
7 -F FORMAT, --list-format=FORMAT
8                       formats --list, choices: short, normal, nested
9 -l, --list            print list of possible commands and exit
10 --set=KEY=VALUE,...  comma separated KEY=VALUE pairs to set Fab env vars
11 --shortlist          alias for -F short --list
12 -V, --version        show program's version number and exit
13 -a, --no_agent       don't use the running SSH agent
14 -A, --forward-agent  forward local agent to remote end
15 --abort-on-prompts   abort instead of prompting (for password, host, etc)
16 -c PATH, --config=PATH
17                       specify location of config file to use
18 -D, --disable-known-hosts
19                       do not load user known_hosts file
20 -f PATH, --fabfile=PATH
21                       python module file to import, e.g. '../other.py'
22 --hide=LEVELS        comma-separated list of output levels to hide
23 -H HOSTS, --hosts=HOSTS
24                       comma-separated list of hosts to operate on
25 -i PATH              path to SSH private key file. May be repeated.
26 -k, --no-keys        don't load private key files from ~/.ssh/
27 --keepalive=N        enables a keepalive every N seconds
28 --linewise           print line-by-line instead of byte-by-byte
29 -n M, --connection-attempts=M
30                       make M attempts to connect before giving up
31 --no-pty             do not use pseudo-terminal in run/sudo
32 -p PASSWORD, --password=PASSWORD
```



```

33                                     password for use with authentication and/or sudo
34 -P, --parallel                       default to parallel execution method
35 --port=PORT                          SSH connection port
36 -r, --reject-unknown-hosts          reject unknown hosts
37                                     reject unknown hosts
38 -R ROLES, --roles=ROLES              comma-separated list of roles to operate on
39                                     comma-separated list of roles to operate on
40 -s SHELL, --shell=SHELL              specify a new shell, defaults to '/bin/bash -l -c'
41                                     specify a new shell, defaults to '/bin/bash -l -c'
42 --show=LEVELS                        comma-separated list of output levels to show
43 --skip-bad-hosts                     skip over hosts that can't be reached
44 --ssh-config-path=PATH               Path to SSH config file
45                                     Path to SSH config file
46 -t N, --timeout=N                   set connection timeout to N seconds
47 -u USER, --user=USER                username to use when connecting to remote hosts
48 -w, --warn-only                      warn, instead of abort, when commands fail
49 -x HOSTS, --exclude-hosts=HOSTS     comma-separated list of hosts to exclude
50                                     comma-separated list of hosts to exclude
51 -z INT, --pool-size=INT              number of concurrent processes to use in parallel mode
52                                     number of concurrent processes to use in parallel mode

```

You can also declare and call functions which take arguments. For example here is a fabfile that allows you to apt-get install packages remotely:

```

1 from fabric.api import env, sudo
2
3 env.hosts = [ '192.168.1.100', '192.168.1.101', '192.168.1.102' ]
4 env.user = "you"
5 env.password = "mysecret"
6 env.parallel = True
7
8 def install(package):
9     sudo('apt-get -y install %s' % package)

```

and you can run it with:

```
1 fab install:python
```

A nice feature of Fabric is the ability to transfer files between the local and the remote host. Here is an example to copy a local file and unzip it at each at the remote destinations:

```

1 from fabric.api import env, sudo
2 from fabric.operations import put
3
4 env.hosts = [ '192.168.1.100', '192.168.1.101', '192.168.1.102' ]
5 env.user = 'you'
6 env.password = 'ysecret'
7 env.parallel = True

```

74 WEB AUTOMATION WITH PYTHON

```
8
9 def deploy(filename):
10     put(filename, '/home/you/destination/path')
11     run('unzip /home/you/destination/path/%s' % filename)
```

which you can use with

```
1 > fab deploy:myfile.zip
```

If you need to change folders:

```
1 from __future__ import with_statement
2 from fabric.api import run
3
4 dir = '/path/to/mydir/'
5
6 with cd(dir):
7     run('ls -l')
```

Fabric can handle errors:

```
1 from __future__ import with_statement
2 from fabric.api import local, settings, abort
3 from fabric.contrib.console import confirm
4
5 def test():
6     with settings(warn_only=True):
7         result = local('./manage.py test my_app', capture=True)
8         if result.failed and not confirm('`Tests failed. Continue anyway?`'):
9             abort('`Aborting at user request.`')
```

8

User Management

To create a new user:

```
1 >>> import os
2 >>> import crypt
3 >>> password = 'testpassword'
4 >>> username = 'someuser'
5 >>> encrypted_password = crypt.crypt(password, '22')
6 >>> os.system('useradd -m -p %s %s' % (encrypted_password, username))
```

(the `-m` option creates the home user if does not exist).

To give superuser permissions to the new user use shell command `visudo` to edit `/etc/sudoers` file:

```
1 > export VISUAL=nano; visudo
2 > visudo
3 # users network=(group) commands
4
5 # joe mike and jude can run any command (ALL) from (ALL) terminals
6 joe, mike, jude ALL=ALL
7
8 # tim and ales can run any command from any terminal as memerbs of the www-data
   group
9 tim, alex ALL=(www-data) ALL
10
11 # jack can can run any command from any group from any machine in
   10.1.2.0/255.255.255.0
12 jack 10.1.2.0/255.255.255.0=(ALL) ALL
13
14 # cathy can print from any machine
15 cathy ALL=/usr/sbin/lpc, /usr/bin/lprm
```

76 WEB AUTOMATION WITH PYTHON

```
16
17 # john can do everything and is not even prompted a password!
18 john ALL=(ALL) NOPASSWD: ALL
```

To automate remote tasks we want to be able to login into a remote host without being prompted a password. We can achieve this by creating a certificate for the client machine. On the client machine:

```
1 > mkdir ~/.ssh
2 > chmod 700 ~/.ssh
3 > ssh-keygen -t rsa -b 4096
```

it will produce the following and as some questions, you can just press enter:

```
1 Generating public/private rsa key pair.
2 Enter file in which to save the key (/home/username/.ssh/id_rsa):
3 Enter passphrase (empty for no passphrase):
4 Enter same passphrase again:
5 Your identification has been saved in /home/username/.ssh/id_rsa.
6 Your public key has been saved in /home/username/.ssh/id_rsa.pub.
```

Then use the following command

```
1 > ssh-copy-id <username>@<host>
```

which will automatically ssh into <username>@<host> and appends the local ~/.ssh/id_rsa.pub to the remote ~/.ssh/authorized_keys.

Troubleshooting. You may need to set permission on the remote host:

```
1 > chmod go-w ~/
2 > chmod 700 ~/.ssh
3 > chmod 600 ~/.ssh/authorized_keys
```

You may need to tell remote host to add the key:

```
1 > ssh-add
```

Now you can:

```
1 > ssh <username>@<host>
```

Always on the remote host, makes sure the file /etc/ssh/sshd_config contains:

```
1 PubkeyAuthentication yes
2 RSAAuthentication yes
```

You may need to restart the daemon. On the host:

```
1 sudo /etc/init.d/sshd restart
```

9

EC2 and OpenStack

A virtual machine (VM) is a software implementation of a machine (i.e. a computer) that executes programs like a physical machine. Virtual machines are separated into two major classifications, based on their use and degree of correspondence to any real machine:

A system virtual machine provides a complete system platform which supports the execution of a complete operating system (OS). These usually emulate an existing architecture, and are built with the purpose of either providing a platform to run programs where the real hardware is not available for use (for example, executing software on otherwise obsolete platforms), or of having multiple instances of virtual machines leading to more efficient use of computing resources, both in terms of energy consumption and cost effectiveness (known as hardware virtualization, the key to a cloud computing environment), or both.

A process virtual machine (also, language virtual machine) is designed to run a single program, which means that it supports a single process. Such virtual machines are usually closely suited to one or more programming languages and built with the purpose of providing program portability and flexibility (amongst other things). An essential characteristic of a virtual machine is that the software running inside is limited to the resources and abstractions provided by the virtual machine—it cannot break out of its virtual environment.

Multiple VMs each running their own operating system (called guest operating system) are frequently used in server consolidation, where different services that used to run on individual machines to avoid interference are instead run in separate VMs on the same physical machine.

Software to run virtual machines include VMWare, VirtualBox, and KVM, and XEN.

There are many web services that allows customers to rent computers on which to run their own computer applications. The biggest is probably Amazon EC2.

9.1 EC2 Commad Line API Tools

Display all available ec2 images

```
ec2-describe-images -x all
```

Generate a keypair which will be instead of a root password

```
ec2-create-keypair mb-keypair
```

Start running a EC2 Instance, this is when billing starts

```
ec2-run-instances ami-a9fe1bc0 -k mbkeypair
```

Open up the port 22 for ssh, similiary 80 for port 80

```
ec2-authorize default -p 22
```

Display list of ec2 instances running

```
ec2-describe-instances
```

Terminate an ec2 running instances. This is where billing ends

```
ec2-terminate-instances ami-a9fe1bc0
```

OpenStack is a free and open software for creating a cloud infrasture like amazon EC2. The technology consists of a series of interrelated projects that control pools of processing, storage, and networking resources throughout a datacenter, all managed through a dashboard that gives administrators

control while empowering its users to provision resources through a web interface.

Nova is a controller for cloud computing platforms, supports OpenStack and EC2.

```
>>> from novaclient import OpenStack
>>> nova = OpenStack(USERNAME, PASSWORD, AUTH_URL)
>>> nova.servers.list()
[<Server: buildslave-ubuntu-9.10>]

>>> nova.flavors.list()
[<Flavor: 256 server>,
 <Flavor: 512 server>,
 <Flavor: 1GB server>,
 <Flavor: 2GB server>,
 <Flavor: 4GB server>,
 <Flavor: 8GB server>,
 <Flavor: 15.5GB server>]
>>> fl = nova.flavors.find(ram=512)
>>> nova.servers.create("my-server", flavor=fl)
<Server: my-server>
```

For more info: <https://github.com/openstack/python-novaclient>

Boto is another python library to manage virtualization. While Nova is designed for OpenStack and supports EC2 via a compatibility layer, Boto is designed for EC2 and supports OpenStack via a compatibility layer.

Both Nova and Boto can run from command line:

```
> nova boot myUbuntuServer --image "3afe97b2-26dc-49c5-a2cc-a2fc8d80c001" --
```


Bibliography

[1] <http://www.python.org>

[2] <http://farmdev.com/talks/unicode/>

[3] <http://www.crummy.com/software/BeautifulSoup/>

[4] <http://wwwsearch.sourceforge.net/mechanize/>

[5] <http://docs.fabfile.org/en/1.8/>

[6] <https://help.ubuntu.com/community/SSH/OpenSSH/Keys>

[7] <http://www.garron.me/en/linux/visudo-command-sudoers-file-sudo-default.html>