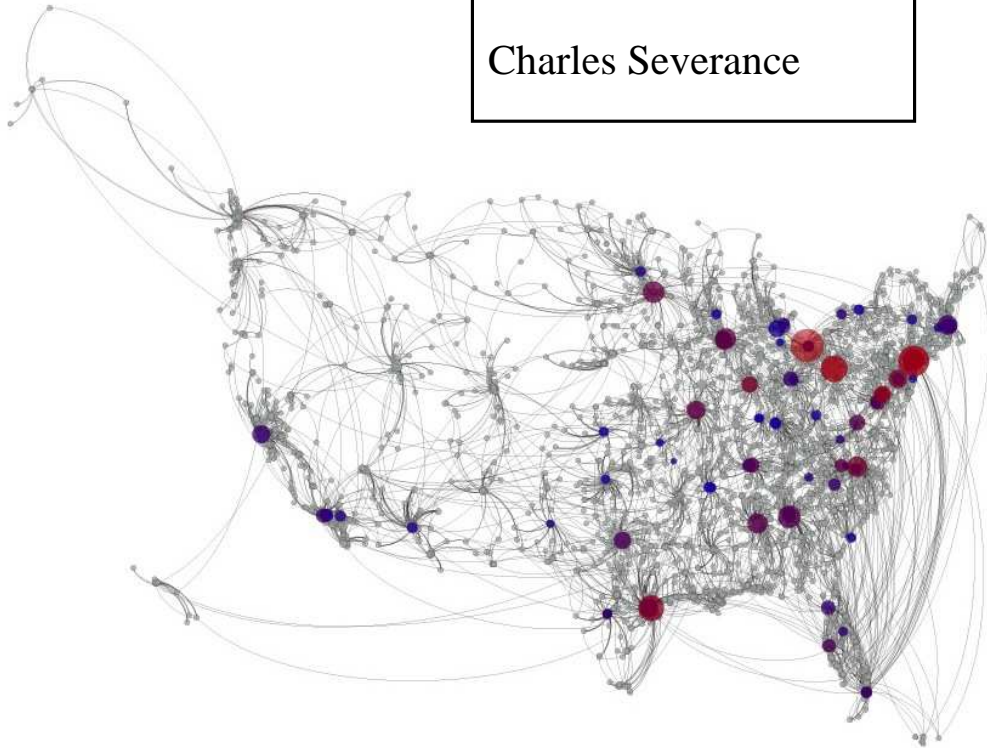


Python for Informatics
Exploring Data

Charles Severance



Python for Informatics

Exploring Information

Version 0.0.3

Charles Severance



Copyright © 2009, 2010 Charles Severance.

Printing history:

December 2009: Begin to produce *Python for Informatics: Exploring Information* by re-mixing *Think Python: How to Think Like a Computer Scientist*

June 2008: Major revision, changed title to *Think Python: How to Think Like a Computer Scientist*.

August 2007: Major revision, changed title to *How to Think Like a (Python) Programmer*.

April 2002: First edition of *How to Think Like a Computer Scientist*.

This work is licensed under a Creative Commons Attribution-Share Alike 3.0 Unported License. This license is available at creativecommons.org/licenses/by-sa/3.0/.

The original form of this book is \LaTeX source code. Compiling this \LaTeX source has the effect of generating a device-independent representation of a textbook, which can be converted to other formats and printed.

The \LaTeX source for the *Think Python: How to Think Like a Computer Scientist* version of this book is available from <http://www.thinkpython.com>.

The \LaTeX source for the *Python for Informatics: Exploring Information* version of the book is available (for the moment) from <http://source.sakaiproject.org/contrib//csev/trunk/pyinf/>.

The cover images were provided by Dr. Lada Adamic and are used with permission.

Preface

Python for Informatics: Remixing an Open Book

It is quite natural for academics who are continuously told to “publish or perish” to want to always create something from scratch that is their own fresh creation. This book is an experiment in not starting from scratch, but instead “re-mixing” the book titled *Think Python: How to Think Like a Computer Scientist* written by Allen B. Downey, Jeff Elkner and others.

In December of 2009, I was preparing to teach **SI502 - Networked Programming** at the University of Michigan for the fifth semester in a row and decided it was time to write a Python textbook that focused on exploring data instead of understanding algorithms and abstractions. My goal in SI502 is to teach people life-long data handling skills using Python. Few of my students were planning to be professional computer programmers. Instead, they planned be librarians, managers, lawyers, biologists, economists, etc. who happened to want to skillfully use technology in their chosen field.

I never seemed to find the perfect data-oriented Python book for my course so I set out to write just such a book. Luckily at a faculty meeting three weeks before I was about to start my new book from scratch over the holiday break, Dr. Atul Prakash showed me the *Think Python* book which he had used to teach his Python course that semester. It is a well-written Computer Science text with a focus on short, direct explanations and ease of learning.

As the copyright holder of *Think Python*, Allen has given me permission to change the book’s license from the GNU Free Documentation License to the more recent Creative Commons Attribution — Share Alike license. This follows a general shift in open documentation licenses moving from the GFDL to the CC-BY-SA (i.e. Wikipedia). Using the CC-BY-SA license maintains the book’s strong copyleft tradition while making it even more straightforward for new authors to reuse this material as they see fit.

I expect that by the time I am done with *Python for Informatics* over fifty percent of the book will be new. The overall structure will be changed to get to doing data analysis problems as quickly as possible and have a series of running examples and exercises about data analysis. Then I will add chapters on regular expressions, data visualization, working with spreadsheet data, structured query language using SQLite, web scraping, and calling REST-based Application Program Interfaces.

The ultimate goal in the shift from a Computer Science to an Informatics focus is to pull topics into the first programming class that can be applied even if one chooses not to become a professional programmer.

What is interesting even with this change of focus is how much of the original *Think Python* book material is directly relevant to this book and how much will fit right into *Python for Informatics* with virtually no change.

By starting with the *Think Python* book, I don't have to write the basic descriptions of the Python language or how to debug programs and instead focus on the topical material that is the value-add of *Python for Informatics*.

Students who find this book interesting and want to further explore a career as a professional programmer should probably look at the *Think Python* book. Because there is a lot of overlap between the two books, you will quickly pick up skills in the additional areas of Computer Science which are covered in *Think Python*. And given that the books have a similar writing style and at times have identical text and examples, you should be able to pick up these new topics with a minimum of effort.

I hope that this book serves an example of why open materials are so important to the future of education, and want to thank Allen B. Downey and Cambridge University Press for their forward looking decision to make the book available under an open Copyright. I hope they are pleased with the results of my efforts and I hope that you the reader are pleased with *our* collective efforts.

Charles Severance
www.dr-chuck.com
December 19, 2009

Charles Severance is a Clinical Assistant Professor at the University of Michigan School of Information.

Draft Version Instructions

The copy of this book you are looking at is currently a draft and still in development. The general roadmap for the rest of the development book is as follows:

- Teach SI502 - Networked Programming at University of Michigan Winter 2010. The first 10 chapters of the book will be used for the first four weeks of the course. At least three more chapters will be written for SI502 and distributed during the semester that line up with the topics in the second half of SI502 (Networked Programming, Databases, and Using Web Services).
- There are four more chapters planned at some point (Advanced Functions, Regular Expressions, Automating Common Tasks, and Visualizing data). These are not currently in the scope of SI502 for Winter 2010.

Like all books being written and used in a course at the same time, student feedback is essential to producing a strong book. So I hope that students will look at the book and

help me find simple errors, places where ideas jump too fast, improvements in the glossary, debugging, and exercises in each chapter.

You can also send comments to csev (at) umich.edu at any time.

Thanks in advance for your patience and assistance.

Preface for “Think Python”

The strange history of “Think Python”

(Allen B. Downey)

In January 1999 I was preparing to teach an introductory programming class in Java. I had taught it three times and I was getting frustrated. The failure rate in the class was too high and, even for students who succeeded, the overall level of achievement was too low.

One of the problems I saw was the books. They were too big, with too much unnecessary detail about Java, and not enough high-level guidance about how to program. And they all suffered from the trap door effect: they would start out easy, proceed gradually, and then somewhere around Chapter 5 the bottom would fall out. The students would get too much new material, too fast, and I would spend the rest of the semester picking up the pieces.

Two weeks before the first day of classes, I decided to write my own book. My goals were:

- Keep it short. It is better for students to read 10 pages than not read 50 pages.
- Be careful with vocabulary. I tried to minimize the jargon and define each term at first use.
- Build gradually. To avoid trap doors, I took the most difficult topics and split them into a series of small steps.
- Focus on programming, not the programming language. I included the minimum useful subset of Java and left out the rest.

I needed a title, so on a whim I chose *How to Think Like a Computer Scientist*.

My first version was rough, but it worked. Students did the reading, and they understood enough that I could spend class time on the hard topics, the interesting topics and (most important) letting the students practice.

I released the book under the GNU Free Documentation License, which allows users to copy, modify, and distribute the book.

What happened next is the cool part. Jeff Elkner, a high school teacher in Virginia, adopted my book and translated it into Python. He sent me a copy of his translation, and I had the unusual experience of learning Python by reading my own book.

Jeff and I revised the book, incorporated a case study by Chris Meyers, and in 2001 we released *How to Think Like a Computer Scientist: Learning with Python*, also under the

GNU Free Documentation License. As Green Tea Press, I published the book and started selling hard copies through Amazon.com and college book stores. Other books from Green Tea Press are available at greenteapress.com.

In 2003 I started teaching at Olin College and I got to teach Python for the first time. The contrast with Java was striking. Students struggled less, learned more, worked on more interesting projects, and generally had a lot more fun.

Over the last five years I have continued to develop the book, correcting errors, improving some of the examples and adding material, especially exercises. In 2008 I started work on a major revision—at the same time, I was contacted by an editor at Cambridge University Press who was interested in publishing the next edition. Good timing!

I hope you enjoy working with this book, and that it helps you learn to program and think, at least a little bit, like a computer scientist.

Acknowledgements for “Think Python”

(Allen B. Downey)

First and most importantly, I thank Jeff Elkner, who translated my Java book into Python, which got this project started and introduced me to what has turned out to be my favorite language.

I also thank Chris Meyers, who contributed several sections to *How to Think Like a Computer Scientist*.

And I thank the Free Software Foundation for developing the GNU Free Documentation License, which helped make my collaboration with Jeff and Chris possible.

I also thank the editors at Lulu who worked on *How to Think Like a Computer Scientist*.

I thank all the students who worked with earlier versions of this book and all the contributors (listed in an Appendix) who sent in corrections and suggestions.

And I thank my wife, Lisa, for her work on this book, and Green Tea Press, and everything else, too.

Allen B. Downey
Needham MA

Allen Downey is an Associate Professor of Computer Science at the Franklin W. Olin College of Engineering.

Contents

Preface	v
1 Why should you learn to write programs?	1
1.1 Creativity and motivation	2
1.2 Computer hardware architecture	3
1.3 Understanding programming	4
1.4 The Python programming language	5
1.5 What is a program?	6
1.6 What is debugging?	7
1.7 Building “sentences” in Python	9
1.8 The first program	11
1.9 Debugging	11
1.10 Glossary	12
1.11 Exercises	13
2 Variables, expressions and statements	15
2.1 Values and types	15
2.2 Variables	16
2.3 Variable names and keywords	17
2.4 Statements	18
2.5 Operators and operands	18
2.6 Expressions	19

2.7	Order of operations	20
2.8	Modulus operator	20
2.9	String operations	21
2.10	Asking the user for input	21
2.11	Comments	22
2.12	Choosing mnemonic variable names	23
2.13	Debugging	24
2.14	Glossary	25
2.15	Exercises	26
3	Conditional execution	29
3.1	Boolean expressions	29
3.2	Logical operators	30
3.3	Conditional execution	30
3.4	Alternative execution	31
3.5	Chained conditionals	31
3.6	Nested conditionals	32
3.7	Catching exceptions using try and except	32
3.8	Short circuit evaluation of logical expressions	34
3.9	Debugging	35
3.10	Glossary	36
3.11	Exercises	37
4	Functions	39
4.1	Function calls	39
4.2	Built-in functions	39
4.3	Type conversion functions	40
4.4	Random numbers	41
4.5	Optional parameters	41
4.6	Glossary	42

Contents	xi
4.7 Math functions	42
4.8 Adding new functions	43
4.9 Definitions and uses	45
4.10 Flow of execution	45
4.11 Parameters and arguments	46
4.12 Fruitful functions and void functions	47
4.13 Why functions?	48
4.14 Debugging	49
4.15 Glossary	49
4.16 Exercises	50
5 Iteration	51
5.1 Updating variables	51
5.2 The while statement	51
5.3 Infinite loops	52
5.4 “Infinite loops” and break	53
5.5 Finishing iterations with continue	55
5.6 Definite loops using for	55
5.7 Loop patterns	56
5.8 Debugging	59
5.9 Glossary	59
5.10 Exercises	60
6 Strings	61
6.1 A string is a sequence	61
6.2 Getting the length of a string using len	62
6.3 Traversal through a string with a for loop	62
6.4 String slices	63
6.5 Strings are immutable	64
6.6 Searching	65

6.7	Looping and counting	65
6.8	The <code>in</code> operator	66
6.9	String comparison	66
6.10	<code>string</code> methods	66
6.11	Parsing strings	69
6.12	Format operator	69
6.13	Debugging	70
6.14	Glossary	73
6.15	Exercises	74
7	Files	77
7.1	Persistence	77
7.2	Opening files	78
7.3	Text files and lines	78
7.4	Reading files	80
7.5	Searching through a file	81
7.6	Letting the user choose the file name	83
7.7	Using <code>try</code> , <code>catch</code> , and <code>open</code>	83
7.8	Writing files	85
7.9	Debugging	86
7.10	Glossary	86
7.11	Exercises	86
8	Lists	89
8.1	A list is a sequence	89
8.2	Lists are mutable	90
8.3	Traversing a list	91
8.4	List operations	91
8.5	List slices	92
8.6	List methods	92

Contents	xiii
8.7 Deleting elements	93
8.8 Lists and strings	94
8.9 Parsing lines	95
8.10 Objects and values	96
8.11 Aliasing	97
8.12 List arguments	97
8.13 Debugging	99
8.14 Glossary	102
8.15 Exercises	103
9 Dictionaries	105
9.1 Dictionary as a set of counters	107
9.2 Dictionaries and files	108
9.3 Looping and dictionaries	109
9.4 Advanced text parsing	110
9.5 Debugging	112
9.6 Glossary	113
9.7 Exercises	113
10 Tuples	115
10.1 Tuples are immutable	115
10.2 Comparing tuples	116
10.3 Tuple assignment	117
10.4 Dictionaries and tuples	119
10.5 Multiple assignment with dictionaries	119
10.6 The most common words	120
10.7 Using tuples as keys in dictionaries	121
10.8 Sequences: strings, lists, and tuples—Oh My!	122
10.9 Debugging	123
10.10 Glossary	124
10.11 Exercises	125

11 Automating common tasks on your computer	127
11.1 File names and paths	127
11.2 Example: Cleaning up a photo directory	128
11.3 Command line arguments	133
11.4 Pipes	135
11.5 Glossary	135
11.6 Exercises	136
12 Networked programs	139
12.1 HyperText Transport Protocol - HTTP	139
12.2 The World's Simplest Web Browser	140
12.3 Retrieving web pages with urllib	141
12.4 Parsing HTML and scraping the web	142
12.5 Glossary	144
12.6 Exercises	145
13 Using Web Services	147
13.1 eXtensible Markup Language - XML	147
13.2 Parsing XML	147
13.3 Looping through nodes	148
13.4 Application Programming Interfaces (API)	149
13.5 Twitter web services	150
13.6 Handling XML data from an API	152
13.7 Glossary	153
13.8 Exercises	154
14 Using databases and Structured Query Language (SQL)	155
14.1 What is a database?	155
14.2 Database concepts	155
14.3 SQLite Database Browser	156

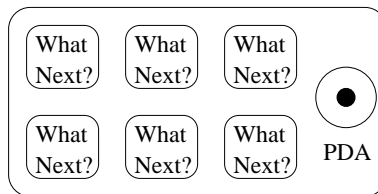
14.4	Creating a database table	156
14.5	Structured Query Language (SQL) summary	159
14.6	Spidering Twitter using a database	160
14.7	Basic data modeling	165
14.8	Programming with multiple tables	166
14.9	Three kinds of keys	171
14.10	Using JOIN to retrieve data	172
14.11	Summary	174
14.12	Debugging	174
14.13	Glossary	175
15	Advanced functions	177
15.1	Return values	177
15.2	Tuples as return values	178
15.3	Variable-length argument tuples	179
15.4	Variables and parameters are local	180
15.5	Global variables	180
15.6	Incremental development	182
15.7	Composition	184
15.8	Stack diagrams	185
15.9	Boolean functions	186
15.10	Optional parameters	186
15.11	Debugging	187
15.12	Glossary	188
15.13	Exercises	189
A	Debugging	191
A.1	Syntax errors	191
A.2	Runtime errors	193
A.3	Semantic errors	196
B	Contributor List	199

Chapter 1

Why should you learn to write programs?

Writing programs (or programming) is a very creative and rewarding activity. You can write programs for many reasons ranging from making your living to solving a difficult data analysis problem to having fun to helping someone else solve a problem. This book assumes that *everyone* needs to know how to program and that once you know how to program, you will figure out what you want to do with your newfound skills.

We are surrounded in our daily lives with computers ranging from laptops to cell phones. We can think of these computers as our “personal assistants” who can take care of many things on our behalf. The hardware in our current-day computers is essentially built to continuously ask us the question, “What would you like me to do next?”.



Programmers add an operating system and a set of applications to the hardware and we end up with a Personal Digital Assistant that is quite helpful and capable of helping many different things.

Our computers are fast and have vast amounts of memory and could be very helpful to us if we only knew the language to speak to explain to the computer what we would like it to “do next”. If we knew this language we could tell the computer to do tasks on our behalf that were repetitive. Interestingly, the kinds of things computers can do best are often the kinds of things that we humans find boring and mind-numbing.

For example, look at the first three paragraphs of this chapter and tell me the most commonly used word and how many times the word is used. While you were able to read and

Thank You for previewing this eBook

You can read the full version of this eBook in different formats:

- HTML (Free /Available to everyone)
- PDF / TXT (Available to V.I.P. members. Free Standard members can access up to 5 PDF/TXT eBooks per month each month)
- Epub & Mobipocket (Exclusive to V.I.P. members)

To download this full book, simply select the format you desire below

