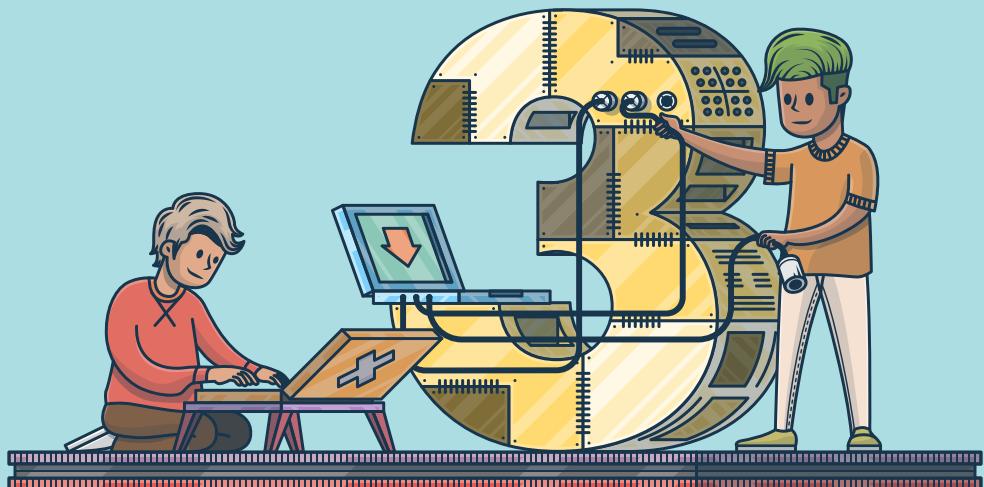


PYTHON BASICS



**A PRACTICAL INTRODUCTION
TO PYTHON 3**

FOURTH EDITION

BY THE REALPYTHON.COM TUTORIAL TEAM
FLETCHER HEISLER, DAVID AMOS, DAN BADER, JOANNA JABLONSKI

Python Basics: A Practical Introduction to Python 3

Real Python

Python Basics

Fletcher Heisler, David Amos, Dan Bader, Joanna Jablonski

Copyright © Real Python (realpython.com), 2012–2020

For online information and ordering of this and other books by Real Python, please visit realpython.com. For more information, please contact us at info@realpython.com.

ISBN: 9781775093329 (paperback)

ISBN: 9781775093336 (electronic)

Cover design by Aldren Santos

“Python” and the Python logos are trademarks or registered trademarks of the Python Software Foundation, used by Real Python with permission from the Foundation.

Thank you for downloading this ebook. This ebook is licensed for your personal enjoyment only. This ebook may not be re-sold or given away to other people. If you would like to share this book with another person, please purchase an additional copy for each recipient. If you’re reading this book and did not purchase it, or it was not purchased for your use only, then please return to realpython.com/pybasics-book and purchase your own copy. Thank you for respecting the hard work behind this book.

Updated 2020-02-25 We would like to thank our early access readers for their excellent feedback: Zoheb Ainapore, Marc, Ricky Mitchell, Meir Guttman, Robert Livingston, Ricky, Jeffrey Hansen, Albrecht, Larry Eisenberg, Kilimandaros, Joanna Jablonski, Mursalin Simpson, Xu Chunyang, Ward Walker, W., Vlad, jima, Vivek, Srinivasan Samuel, Patrick Starrenburg, marp, Jorge Alberch, Edythe, Miguel Galán, Tom Carnevale, Florent, Albrecht Kadauke, Hans van Nielen, Youri Torchalski, Gavin, Karen H Calhoun MD, Roman, Robert Robb Livingston, Terence Phillips, Nico, Daniel, Kumaran Rajendhiran, Ty Wait, david fullerton, Robert, Nicklas, Jacob Andersen, Mario,

Alejandro Ramos, Beni_begin, AJ, Melvin, Sean Yang, Sean, Velu.V, Peter Cavallaro, Charlie Browning 3, Milind Mahajani, Jason Barnes, Lucien Boland, Adam bretel, William, Veltaine, Jerry Petrey, James, Raymond E Rogers, ty wait, Bimperng Uen, CJ Hwang, Guido, Evan, Dave, Miguel Galan, Han Qi, Jim Bremner, Matt Chang, Daniel Drazan, Cole, Bob, Reed Howald, Edward Duarte, Mike Parker, Aart Kleinendorst, rock, Johnnny, Rock Lee, Dusan Ranisavljev, Grant, Jack, Reinhard, Ceejay Cervantes, David, Vivek Vashist, Mark, Dan, Garrett, Peter, Jun Lee, James Silk, Nik Singhal, Charles, Allard Schmidt, Jeff Desalle, Miguel, Steve Poe, Jonathan Seubert, Marc Poulin, MELVIN, Idris, Lucas, John Chirico, Wynette Espinosa, J.P., Gregory, Mark Edgeller, David Melanson, Raul Pena, Darrell, Shriram, Tom Flynn, Velu, michael lindsey, Sulo Kolehmainen, Michael, Jay, Richard, Milos "Ozzyx" Kosik, hans de Cocq, Glen Mules, Nathan Lundner, Phil, Shubh, Puwei Wang, Alex Mück, Alex, Hitoshi, Bruno F. De Lima, Dario David, Rajesh, Haroldas Valčiukas, GVeltaine, Susan Fowle, Jared Simms, Nathan Collins, Dylan, Les Churchman (luckyles in the Pythonistacafe), Stephane LI-THIAO-TE, Frank P, Paul, Damien Murtagh, Jason, Thắng Lê Quang, Neill, Lele, charles wilson, Damien, Christian, Jon, Andreas Kreisig, Marco, Mario Panagiotopoulos, nerino, Mariusz, Thomas, Mihhail, Mikönig, Fabio, Scott, Pedro Torres, Mathias Johansson, Joshua S., Mathias, scott, David Koppy, Rohit Bharti, Phillip Douglas, John Stephenson, Jeff Jones, George Mast, Allards, Palak, Nikola N., Palak Kalsi, Annekathrin, Tsung-Ju Yang, Nick Huntington, Sai, Jordan, Wim Alsemgeest, DJ, Bob Harris, Martin, Andrew, Reggie Smith, Steve Santy, tstalin22@gmail.com, Mohee Jarada, Mark Arzaga, Poulose Matthen, Brent Gordon, Gary Butler, Bryant, Dana, Koajck, Reggie, Luis Bravo, Elijah, Nikolay, Eric Lietsch, Fred Janssen, Don Stillwell, Gaurav Sharma, Mike, Mike McKenna, karthik babu, bulat, Bulat Mansurov, August Trillanes, Darren Saw, Jagadish, Nathan Eger, Kyle, Tejas Shetty, Baba Sariffodeen, Don, Ian, Ian Barbour, Redhouane, Wayne Rosing, Emanuel, Toigongonbai, Jason Castillo, krishna chaitanya swamy kesavarapu, Corey Huguley, Nick, w.g.sneddon@gmail.com, xuchunyang, Daniel BUIS, kenneth, Leodanis Pozo Ramos, John Phenix, Linda Moran, W Laleau, Troy Flynn, Heber Nielsen, ROCK, Mike LeRoy, Thomas Davis, Jacob, Szabolcs

Sinka, kalaiselvan, Leanne Kuss, Andrey, omar, Jason Woden, David Cebalo, John Miller, David Bui (newbie), Nico Zanferrari, Ariel, Boris, Boris Ender, Charlie3, Ossy, Matthias Kuehl, Scott Koch, Jesus Avina, charlie3, Awadhesh, Andie, Chris Johnson, Malan, Ciro, Thamizhselvan, Neha, Christian Langpap, Ivan, Dr. Craig Levy, H B Robinson, Stéphane, Steve McIlree, Yves, Teresa, Allard, tom cone jr, Dirk, Joachim van der Weijden, Jim Woodward, Christoph Lipka, John Vergelli, Gerry, Lu, Robert R., Vlad, Richard Heatwole, Gabriel, Krzysztof Surowiecki, Alexandra Davis, Jason Voll, and Dwayne Dever. Thank you all!

This is an Early Access version of “Python Basics: A Practical Introduction to Python 3”

With your help we can make this book even better:

At the end of each section of the book you’ll find a “magical” feedback link. Clicking the link takes you to an **online feedback form** where you can share your thoughts with us.

We welcome any and all feedback or suggestions for improvement you may have.

Please feel free to be as terse or detailed as you see fit. All feedback is stored anonymously, but you can choose to leave your name and contact information so we can follow up or mention you on our “Thank You” page.

We use a different feedback link for each section, so we’ll always know which part of the book your notes refer to.

Thank you for helping us make this book an even more valuable learning resource for the Python community.

— Dan Bader, Editor-in-Chief at Real Python

What Pythonistas Say About *Python Basics: A Practical Introduction to Python 3*

“I love [the book]! The wording is casual, easy to understand, and makes the information flow well. I never feel lost in the material, and it’s not too dense so it’s easy for me to review older chapters over and over.

I’ve looked at over 10 different Python tutorials/books/online courses, and I’ve probably learned the most from Real Python!”

— **Thomas Wong**

“Three years later and I still return to my Real Python books when I need a quick refresher on usage of vital Python commands.”

— **Rob Fowler**

“I floundered for a long time trying to teach myself. I slogged through dozens of incomplete online tutorials. I snoozed through hours of boring screencasts. I gave up on countless crusty books from big-time publishers. And then I found Real Python.

The easy-to-follow, step-by-step instructions break the big concepts down into bite-sized chunks written in plain English. The authors never forget their audience and are consistently thorough and detailed in their explanations. I’m up and running now, but I constantly refer to the material for guidance.”

— **Jared Nielsen**

“I love the book because at the end of each particular lesson there are real world and interesting challenges. I just built a savings estimator that actually reflects my savings account – neat!”

— **Drew Prescott**

“As a practice of what you taught I started building simple scripts for people on my team to help them in their everyday duties. When my managers noticed that, I was offered a new position as a developer.

I know there is heaps of things to learn and there will be huge challenges, but I finally started doing what I really came to like.

Once again: MANY THANKS!”

— **Kamil**

“What I found great about the Real Python courses compared to others is how they explain things in the simplest way possible.

A lot of courses, in any discipline really, require the learning of a lot of jargon when in fact what is being taught could be taught quickly and succinctly without too much of it. The courses do a very good job of keeping the examples interesting.”

— **Stephen Grady**

“After reading the first Real Python course I wrote a script to automate a mundane task at work. What used to take me three to five hours now takes less than ten minutes!”

— **Brandon Youngdale**

“Honestly, throughout this whole process what I found was just me looking really hard for things that could maybe be added or improved, but this tutorial is amazing! You do a wonderful job of explaining and teaching Python in a way that people like me, a complete novice, could really grasp.

The flow of the lessons works perfectly throughout. The exercises truly helped along the way and you feel very accomplished when you finish up the book. I think you have a gift for making Python seem more attainable to people outside the programming world.

This is something I never thought I would be doing or learning and with a little push from you I am learning it and I can see that it will be nothing but beneficial to me in the future!”

— **Shea Klusewicz**

“The authors of the courses have NOT forgotten what it is like to be a beginner – something that many authors do – and assume nothing about their readers, which makes the courses fantastic reads. The courses are also accompanied by some great videos as well as plenty of references for extra learning, homework assignments and example code that you can experiment with and extend.

I really liked that there was always full code examples and each line of code had good comments so you can see what is doing what.

I now have a number of books on Python and the Real Python ones are the only ones I have actually finished cover to cover, and they are hands down the best on the market. If like me, you’re not a programmer (I work in online marketing) you’ll find these courses to be like a mentor due to the clear, fluff-free explanations! Highly recommended!”

— **Craig Addyman**

About the Authors

At [Real Python](#) you'll learn real-world programming skills from a community of professional Pythonistas from all around the world.

The [realpython.com](#) website launched in 2012 and currently helps more than two million Python developers each month with free programming tutorials and in-depth learning resources.

Everyone who worked on this book is a *practitioner* with several years of professional experience in the software industry. Here are the members of the Real Python Tutorial Team who worked on *Python Basics*:

Fletcher Heisler is the founder of Hunter2, where he teaches developers how to hack and secure modern web apps. As one of the founding members of Real Python, Fletcher wrote the original version of this book in 2012.

David Amos is a mathematician by training, a data scientist/Python developer by profession, and a coffee junkie by choice. He is a member of the Real Python tutorial team and rewrote large parts of this book to update it to Python 3.

Dan Bader is the owner and Editor in Chief of Real Python and a complete Python nut. When he's not busy working on the Real Python learning platform he helps Python developers take their coding skills to the next level with tutorials, books, and online training.

Joanna Jablonski is the Executive Editor of Real Python. She loves natural languages just as much as she loves programming languages. When she's not producing educational materials to help Python developers level up, she's finding new ways to optimize various aspects of her life.

Contents

Contents	10
Foreword	15
1 Introduction	22
1.1 Why This Book?	23
1.2 About Real Python	25
1.3 How to Use This Book	25
1.4 Bonus Material & Learning Resources	27
2 Setting Up Python	30
2.1 A Note On Python Versions	31
2.2 Windows	32
2.3 macOS	35
2.4 Ubuntu Linux	39
3 Your First Python Program	43
3.1 Write a Python Script	43
3.2 Mess Things Up	49
3.3 Create a Variable	52
3.4 Inspect Values in the Interactive Window	57
3.5 Leave Yourself Helpful Notes	60
3.6 Summary and Additional Resources	63
4 Strings and String Methods	65
4.1 What is a String?	66
4.2 Concatenation, Indexing, and Slicing	73

4.3	Manipulate Strings With Methods	81
4.4	Interact With User Input	88
4.5	Challenge: Pick Apart Your User’s Input	90
4.6	Working With Strings and Numbers	91
4.7	Streamline Your Print Statements	97
4.8	Find a String in a String	99
4.9	Challenge: Turn Your User Into a L33t H4xor	102
4.10	Summary and Additional Resources	103
5	Numbers and Math	105
5.1	Integers and Floating-Point Numbers	106
5.2	Arithmetic Operators and Expressions	110
5.3	Challenge: Perform Calculations on User Input	118
5.4	Make Python Lie to You	119
5.5	Math Functions and Number Methods	121
5.6	Print Numbers in Style	126
5.7	Complex Numbers	129
5.8	Summary and Additional Resources	133
6	Functions and Loops	135
6.1	What is a Function, Really?	136
6.2	Write Your Own Functions	140
6.3	Challenge: Convert Temperatures	149
6.4	Run in Circles	150
6.5	Challenge: Track Your Investments	160
6.6	Understand Scope in Python	161
6.7	Summary and Additional Resources	166
7	Finding and Fixing Code Bugs	168
7.1	Use the Debug Control Window	169
7.2	Squash Some Bugs	176
7.3	Summary and Additional Resources	185
8	Conditional Logic and Control Flow	186
8.1	Compare Values	187
8.2	Add Some Logic	190
8.3	Control the Flow of Your Program	198

8.4	Challenge: Find the Factors of a Number	210
8.5	Break Out of the Pattern	211
8.6	Recover From Errors	215
8.7	Simulate Events and Calculate Probabilities	221
8.8	Challenge: Simulate a Coin Toss Experiment	227
8.9	Challenge: Simulate an Election	227
8.10	Summary and Additional Resources	228
9	Tuples, Lists, and Dictionaries	230
9.1	Tuples Are Immutable Sequences	231
9.2	Lists Are Mutable Sequences	241
9.3	Nesting, Copying, and Sorting Tuples and Lists . . .	254
9.4	Challenge: List of lists	260
9.5	Challenge: Wax Poetic	261
9.6	Store Relationships in Dictionaries	263
9.7	Challenge: Capital City Loop	274
9.8	How to Pick a Data Structure	275
9.9	Challenge: Cats With Hats	276
9.10	Summary and Additional Resources	277
10	Object-Oriented Programming (OOP)	279
10.1	Define a Class	280
10.2	Instantiate an Object	284
10.3	Inherit From Other Classes	291
10.4	Challenge: Model a Farm	301
10.5	Summary and Additional Resources	302
11	Modules and Packages	304
11.1	Working With Modules	305
11.2	Working With Packages	315
11.3	Summary and Additional Resources	326
12	File Input and Output	328
12.1	Files and the File System	329
12.2	Working With File Paths in Python	333
12.3	Common File System Operations	341
12.4	Challenge: Move All Image Files To a New Directory .	358

12.5	Reading and Writing Files	359
12.6	Read and Write CSV Data	374
12.7	Challenge: Create a High Scores List	385
12.8	Summary and Additional Resources	386
13	Installing Packages With Pip	388
13.1	Installing Third-Party Packages With Pip	389
13.2	The Pitfalls of Third-Party Packages	400
13.3	Summary and Additional Resources	401
14	Creating and Modifying PDF Files	403
14.1	Extract Text From a PDF	404
14.2	Extract Pages From a PDF	411
14.3	Challenge: PdfFileSplitter Class	418
14.4	Concatenating and Merging PDFs	419
14.5	Rotating and Cropping PDF Pages	426
14.6	Encrypting and Decrypting PDFs	438
14.7	Challenge: Unscramble A PDF	442
14.8	Create a PDF File From Scratch	442
14.9	Summary and Additional Resources	449
15	Working With Databases	451
15.1	An Introduction to SQLite	452
15.2	Libraries for Working With Other SQL Databases	464
15.3	Summary and Additional Resources	465
16	Interacting With the Web	467
16.1	Scrape and Parse Text From Websites	468
16.2	Use an HTML Parser to Scrape Websites	477
16.3	Interact With HTML Forms	482
16.4	Interact With Websites in Real-Time	489
16.5	Summary and Additional Resources	493
17	Scientific Computing and Graphing	495
17.1	Use NumPy for Matrix Manipulation	496
17.2	Use matplotlib for Plotting Graphs	507
17.3	Summary and Additional Resources	533

18 Graphical User Interfaces	535
18.1 Add GUI Elements With EasyGUI	536
18.2 Example App: PDF Page Rotator	548
18.3 Challenge: PDF Page Extraction Application	555
18.4 Introduction to Tkinter	556
18.5 Working With Widgets	560
18.6 Controlling Layout With Geometry Managers	588
18.7 Making Your Applications Interactive	607
18.8 Example App: Temperature Converter	617
18.9 Example App: Text Editor	622
18.10 Challenge: Return of the Poet	631
18.11 Summary and Additional Resources	633
19 Final Thoughts and Next Steps	635
19.1 Free Weekly Tips for Python Developers	636
19.2 Python Tricks: The Book	636
19.3 Real Python Video Course Library	637
19.4 PythonistaCafe: A Community for Python Developers	638
19.5 Acknowledgements	640

Foreword

Hello and welcome to **Python Basics: A Practical Introduction to Python 3**. I hope you are ready to learn why so many professional and hobbyist developers are drawn to Python and how you can begin using it on your projects, small and large, right away.

This book is targeted at beginners who either know a little programming but not the Python language and ecosystem, as well as complete beginners.

If you don't have a Computer Science degree, don't worry. Fletcher, David, Dan, and Joanna will guide you through the important computing concepts while teaching you the Python basics, and just as importantly, skipping the unnecessary details at first.

Python Is a Full-Spectrum Language

When learning a new programming language, you don't yet have the experience to judge how well it will serve you in the long run. If you are considering Python, let me assure you that this is a good choice. One key reason is that Python is a **full-spectrum** language.

What do I mean by this? Some languages are very good for beginners. They hold your hand and make programming super easy. We can go to the extreme and look at visual languages such as Scratch.

Here you get blocks that represent programming concepts, like variables, loops, method calls, and so on, and you drag and drop them on a visual surface. Scratch may be easy to get started with for sim-

ple programs. But you cannot build professional applications with it. Name one Fortune 500 company that powers its core business logic with Scratch.

Came up empty? Me too—because that would be insanity.

Other languages are incredibly powerful for expert developers. The most popular one in this category is likely C++ and its close relative C. Whatever web browser you used today was likely written in C or C++. Your operating system running that browser was also very likely built with C/C++. Your favorite first-person shooter or strategy video game? You nailed it: C/C++.

You can do amazing things with these languages. But they are wholly unwelcoming to newcomers looking for a gentle introduction.

You might not have read a lot of C++ code. It can almost make your eyes burn. Here's an example, a real albeit complex one:

```
template <typename T>
_Defer<void(*(PID<T>, void (T::*)(void)))>
    (const PID<T>&, void (T::*)(void))>
defer(const PID<T>& pid, void (T::*method)(void))
{
    void (*dispatch)(const PID<T>&, void (T::*)(void)) =
        &process::template dispatch<T>;
    return std::bind(dispatch, pid, method);
}
```

Please, just no.

Both Scratch and C++ are decidedly not what I would call full-spectrum languages. In the Scratch level, it's easy to start but you have to switch to a "real" language to build real applications. Conversely, you can build real apps with C++, yet there is no gentle on-ramp. You dive head first into all the complexity of that language which exists to support these rich applications.

Python, on the other hand, is special. It is a full-spectrum language. We often judge the simplicity of a language based on the “hello world” test. That is, what syntax and actions are necessary to get that language to output “hello world” to the user? In Python, it couldn’t be simpler:

```
print("Hello world")
```

That’s it. However, I find this an unsatisfying test.

The “hello world” test is useful but really not enough to show the power or complexity of a language. Let’s try another example. Not everything here needs to make total sense, just follow along to get the Zen of it. The book covers these concepts and more as you go through. The next example is certainly something you could write near the end.

Here’s the new test: What would it take to write a program that accesses an external website, downloads the content to your app in memory, then displays a subsection of that content to the user? Let’s try that experiment with Python 3 with the help of the `requests` package (which needs to be installed—more on that in chapter 12):

```
import requests
resp = requests.get("https://realpython.com")
html = resp.text
print(html[205:294])
```

Incredibly, that’s it. When run, the output is (something like):

```
<title>Python Tutorials - Real Python</title>
<meta name="author" content="Real Python">
```

This is the easy, getting started side of the spectrum of Python. A few trivial lines and incredible power is unleashed. Because Python has access to so many powerful but well-packaged libraries, such as `requests`, it is often described as *having batteries included*.

So there you have a simple powerful starter example. On the real apps

side of things, we have many incredible applications written in Python as well.

YouTube, the world's most popular video streaming site, is written in Python and processes more than 1,000,000 requests per second. Instagram is another example of a Python application. More close to home, we even have realpython.com and my sites such as talkpython.fm.

This full-spectrum aspect of Python means you can start easy and adopt more advanced features as you need them when your application demands grow.

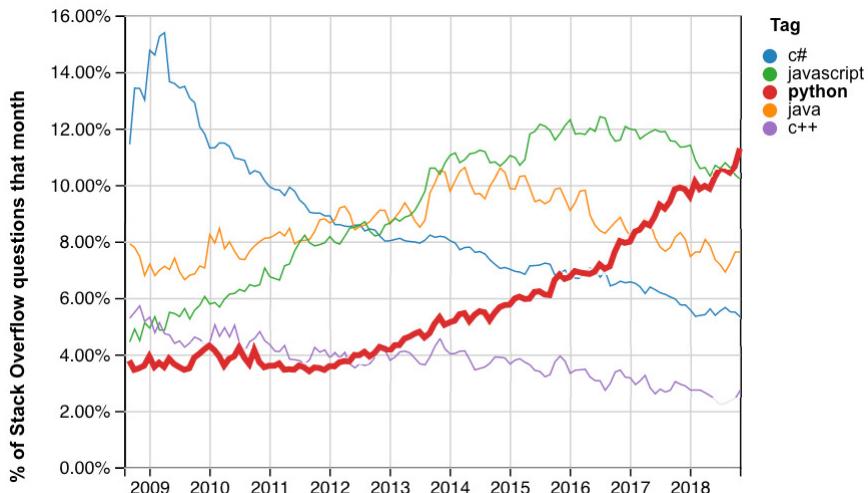
Python Is Popular

You might have heard that Python is popular. On one hand, it may seem that it doesn't really matter how popular a language is if you can build the app you want to build with it.

For better or worse, in software development popularity is a strong indicator of the quality of libraries you will have available as well the number of job openings there are. In short, you should tend to gravitate towards more popular technologies as there will be more choices and integrations available.

So, is Python actually that popular? Yes it is. You'll of course find a lot of hype and hyperbole. But there are plenty of stats to back this one. Let's look at some analytics available and presented by StackOverflow.com.

They run a site called **StackOverflow Trends**. Here you can look at the trends for various technologies by tag. When we compare Python to the other likely candidates you could pick to learn programming, you'll see one is unlike the others:



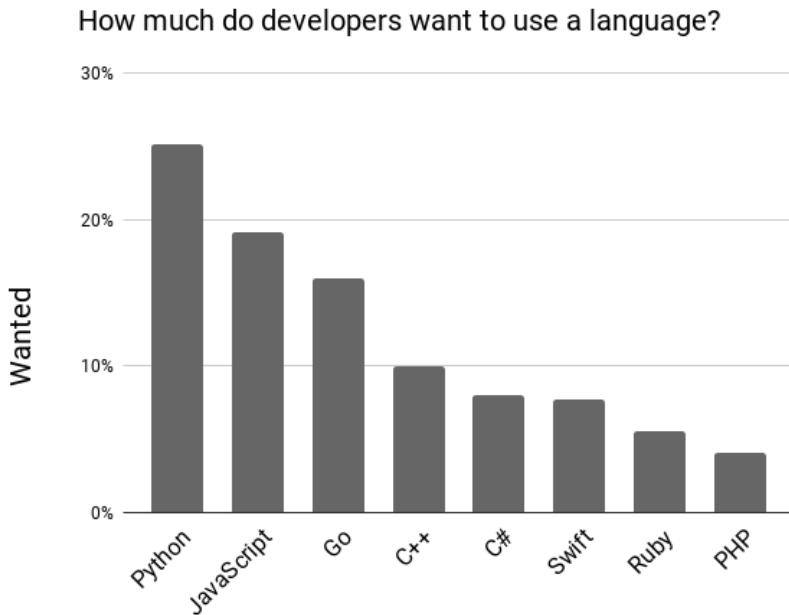
You can explore this chart and create similar charts to this one over at insights.stackoverflow.com/trends.

Notice the incredible growth of Python compared to the flatline or even downward trend of the other usual candidates! If you are betting your future on the success of a given technology, which one would you choose from this list?

That's just one chart, what does it really tell us? Well, let's look at another. StackOverflow does a yearly survey of developers. It's comprehensive and very well done. You can find the full 2018 results at insights.stackoverflow.com/survey/2018/. From that writeup, I'd like to call your attention to a section entitled **Most Loved, Dreaded, and Wanted Languages**. In the most wanted section, you'll find responses for:

Developers who are not developing with the language or technology but have expressed interest in developing with it.

Again, in the graph below, you'll see that Python is topping the charts and well above even second place:



So if you agree with me that the relative popularity of a programming language matters. Python is clearly a good choice.

We Don't Need You to Be a Computer Scientist

One other point I do want to emphasize as you start this journey of learning Python is that we don't need you to be a computer scientist. If that's your goal, great. Learning Python is a powerful step in that direction. But learning programming is often framed in the shape of “we have all these developer jobs going unfilled, we need software developers!”

That may or may not be true. But more importantly for you, programming (even a little programming) can be a superpower for you personally.

To illustrate this idea, suppose you are a biologist. Should you drop out of biology and get a front-end web developer job? Probably not. But having skills such as the one I opened this foreword with, using requests to get data from the web, will be incredible powerful for you as you do biology.

Rather than manually exporting and scraping data from the web or spreadsheets, with Python you can scrape 1,000's of data sources or spreadsheets in the time it takes you to do just one manually. Python skills can be what takes your *biology power* and amplifies it well beyond your colleagues' and makes it your *superpower*.

Dan and Real Python

Finally, let me leave you with a comment on your authors. Dan Bader along with the other Real Python authors work day in and out to bring clear and powerful explanations of Python concepts to all of us via realpython.com.

They have a unique view into the Python ecosystem and are keyed into what beginners need to know.

I'm confident leaving you in their hands on this Python journey. Go forth and learn this amazing language using this great book. Most importantly, remember to have fun!

— **Michael Kennedy**, Founder of Talk Python ([@mkennedy](https://twitter.com/mkennedy))

Chapter 1

Introduction

Welcome to Real Python’s *Python Basics* book, fully updated for Python 3.8! In this book you’ll learn real-world Python programming techniques, illustrated with useful and interesting examples.

Whether you’re new to programming or a professional software developer looking to dive into a new language, this book will teach you all of the practical Python that you need to get started on projects on your own.

No matter what your ultimate goals may be, if you work with a computer at all, you will soon be finding endless ways to improve your life by automating tasks and solving problems through Python programs that you create.

But what’s so great about Python as a programming language? Python is open-source freeware, meaning you can download it for free and use it for any purpose, commercial or not.

Python also has an amazing community that has built a number of additional useful tools you can use in your own programs. Need to work with PDF documents? There’s a comprehensive tool for that. Want to collect data from web pages? No need to start from scratch!

Python was built to be easier to use than other programming lan-

guages. It's usually much easier to read Python code and much faster to write code in Python than in other languages.

For instance, here's some simple code written in C, another commonly used programming language:

```
#include <stdio.h>

int main(void)
{
    printf("Hello, world\n");
}
```

All the program does is show the text `Hello, world` on the screen. That was a lot of work to output one phrase! Here's the same program, written in Python:

```
print("Hello, world")
```

That's pretty simple, right? The Python code is faster to write and easier to read. We find that it looks friendlier and more approachable, too!

At the same time, Python has all the functionality of other languages and more. You might be surprised how many professional products are built on Python code: Instagram, YouTube, Reddit, Spotify, to name just a few.

Not only is Python a friendly and fun language to learn—it also powers the technology behind multiple world-class companies and offers fantastic career opportunities for any programmer who masters it.

1.1 Why This Book?

Let's face it, there's an overwhelming amount of information about Python on the internet.

But many beginners who are studying on their own have trouble fig-

uring out *what* to learn and *in what order* to learn it.

You may be asking yourself, “What should I learn about Python in the beginning to get a strong foundation?” If so, this book is for you—whether you’re a complete beginner or already dabbled in Python or other languages before.

Python Basics is written in plain English and breaks down the core concepts you really need to know into bite-sized chunks. This means you’ll know “enough to be dangerous” with Python, fast.

Instead of just going through a boring list of language features, you’ll see exactly how the different building blocks fit together and what’s involved in building real applications and scripts with Python.

Step by step you’ll master fundamental Python concepts that will help you get started on your journey to learn Python.

Many programming books try to cover every last possible variation of every command which makes it easy for readers to get lost in the details. This approach is great if you’re looking for a reference manual, but it’s a horrible way to learn a programming language. Not only do you spend most of your time cramming things into your head you’ll never use, it also isn’t any fun!

This book is built on the 80/20 principle. We will cover the commands and techniques used in the vast majority of cases and focus on how to program real-world solutions to problems that will help make your life easier.

This way, we guarantee that you will:

- Learn useful programming techniques quickly
- Spend less time struggling with unimportant complications
- Find more practical uses for Python in your own life
- Have more fun in the process

Once you’ve mastered the material in this book, you will have gained

a strong enough foundation that venturing out into more advanced territory on your own will be a breeze.

So dive in! Learn to program in a widely used, free language that can do more than you ever thought was possible.

1.2 About Real Python

At [Real Python](#), you'll learn real-world programming skills from a community of professional Pythonistas from all around the world.

The [realpython.com](#) website launched in 2012 and currently helps more than two million Python developers each month with books, programming tutorials, and other in-depth learning resources.

Everyone who worked on this book is a *Python practitioner* recruited from the Real Python team with several years of professional experience in the software industry.

Here's where you can find Real Python on the web:

- [realpython.com](#)
- [@realpython](#) on Twitter
- The Real Python Email Newsletter

1.3 How to Use This Book

The first half of this book is a quick but thorough overview of all the Python fundamentals. You do not need any prior experience with programming to get started. The second half is focused on finding practical solutions to interesting, real-world coding problems.

As a beginner, we recommend that you go through the first half of this book from start to end. The second half covers topics that don't overlap as much so you can jump around more easily, but the chapters do increase in difficulty as you go along.

If you are a more experienced programmer, then you may find yourself heading toward the second part of the book right away. But don't neglect getting a strong foundation in the basics first and be sure to fill in any knowledge gaps along the way.

Most sections within a chapter are followed by **review exercises** to help you make sure that you've mastered all the topics covered. There are also a number of **code challenges**, which are more involved and usually require you to tie together a number of different concepts from previous chapters.

The practice files that accompany this book also include full solutions to the challenges as well as some of the trickier exercises. But to get the most out of the material, you should try your best to solve the challenge problems on your own before looking at the example solutions.

If you're completely new to programming, you may want to supplement the first few chapters with additional practice. We recommend working through the *Python Fundamentals* tutorials available for free at realpython.com to make sure you are on solid footing.

If you have any questions or feedback about the book, you're always welcome to [contact us](#) directly.

Learning by Doing

This book is all about learning by doing, so be sure to *actually type in* the code snippets you encounter in the book. For best results, we recommend that you avoid copying and pasting the code examples.

You will learn the concepts better and pick up the syntax faster if you type out each line of code yourself. Plus, if you screw up—which is totally normal and happens to all developers on a daily basis—the simple act of correcting typos will help you learn how to debug your code.

Try to complete the review exercises and code challenges on your own before getting help from outside resources. With enough practice, you will master this material—and have fun along the way!

How Long Will It Take to Finish This Book?

If you're already familiar with a programming language you could finish the book in as little as 35 to 40 hours. If you're new to programming you may need to spend up to 100 hours or more. Take your time and don't feel like you have to rush. Programming is a super rewarding, but complex skill to learn. Good luck on your Python journey, we're rooting for you!

1.4 Bonus Material & Learning Resources

Online Resources

This book comes with a number of free bonus resources that you can access at realpython.com/python-basics/resources. On this web page you can also find an errata list with corrections maintained by the Real Python team.

Interactive Quizzes

Most chapters in this book come with a free online quiz to check your learning progress. You can access the quizzes using the links provided at the end of the chapter. The quizzes are hosted on the Real Python website and can be viewed on your phone or computer.

Each quiz takes you through a series of questions related to a particular chapter in the book. Some of them are multiple choice, some will ask you to type in an answer, and some will require you to write actual Python code. As you make your way through each quiz, it keeps score of which questions you answered correctly.

At the end of the quiz you receive a grade based on your result. If you don't score 100% on your first try—don't fret! These quizzes are meant to challenge you and it's expected that you go through them several times, improving your score with each run.

Exercises Code Repository

This book has an accompanying [code repository on the web](#) containing example source code as well as the answers to exercises and code challenges. The repository is broken up by chapter so you can check your code against the solutions provided by us after you finish each chapter. Here's the link:

realpython.com/python-basics/exercises

Example Code License

The example Python scripts associated with this book are licensed under a [Creative Commons Public Domain \(CC0\) License](#). This means that you're welcome to use any portion of the code for any purpose in your own programs.

Note

The code found in this book has been tested with Python 3.8 on Windows, macOS, and Linux.

Formatting Conventions

Code blocks will be used to present example code:

```
# This is Python code:  
print("Hello world!")
```

Terminal commands follow the Unix format:

```
$ # This is a terminal command:  
$ python hello-world.py
```

(Dollar signs are not part of the command.)

Italic text will be used to denote a file name: *hello-world.py*.

Bold text will be used to denote a new or important term.

Keyboard shortcuts will be formatted as follows: **Ctrl + S**.

Menu shortcuts will be formatted as follows: **File > New File**

Notes and Warning boxes appear as follows:

Note

This is a note filled in with placeholder text. The quick brown fox jumps over the lazy dog. The quick brown Python slithers over the lazy hog.

Important

This is a warning also filled in with placeholder text. The quick brown fox jumps over the lazy dog. The quick brown Python slithers over the lazy hog.

Feedback & Errata

We welcome ideas, suggestions, feedback, and the occasional rant. Did you find a topic confusing? Did you find an error in the text or code? Did we leave out a topic you would love to know more about?

We're always looking to improve our teaching materials. Whatever the reason, please send in your feedback at the link below:

realpython.com/python-basics/feedback

[Leave feedback on this section »](#)

Chapter 2

Setting Up Python

This book is about programming computers with Python. You could read this book cover-to-cover and absorb the information without ever touching a keyboard, but you'd miss out on the fun part—coding.

To get the most out of this book, you need to have a computer with Python installed on it and a way to create, edit, and save Python code files.

In this chapter, you will learn how to:

- Install the latest version of Python 3 on your computer
- Open **IDLE**, Python's built-in **I**ntegrated **Development and **L**earning **E**nvironment**

Let's get started!

[Leave feedback on this section »](#)

2.1 A Note On Python Versions

Many operating systems, such as macOS and Linux, come with Python pre-installed. The version of Python that comes with your operating system is called your **system Python**.

The system Python is almost always out-of-date and may not even be a full Python installation. It's essential that you have the most recent version of Python so that you can follow along successfully with the examples in this book.

It's possible to have multiple versions of Python installed on your computer. In this chapter, you'll install the latest version of Python 3 alongside any system Python that may already exist on your machine.

Note

Even if you already have Python 3.8 installed, it is still a good idea to skim this chapter to double check that your environment is set-up for following along with this book.

This chapter is split into three sections: Windows, macOS, and Ubuntu Linux. Find the section for your operating system and follow the steps to get set-up, then skip ahead to the next chapter.

If you have a different operating system, check out Real Python's [Python 3 Installation & Setup Guide](#) to see if your OS is covered.

[Leave feedback on this section »](#)

2.2 Windows

Follow these steps to install Python 3 and open IDLE on Windows.

Important

The code in this book is only tested against Python installed as described in this section.

Be aware that if you have installed Python through some other means, such as Anaconda Python, you may encounter problems when running the code examples.

Install Python

Windows systems do not typically ship with Python pre-installed. Fortunately, installation does not involve much more than downloading the Python installer from the [python.org website](https://www.python.org/) and running it.

Step 1: Download the Python 3 Installer

Open a browser window and navigate to the [download page for Windows at python.org](https://www.python.org/).

Underneath the heading at the top that says *Python Releases for Windows*, click on the link for the *Latest Python 3 Release - Python 3.x.x*. As of this writing, the latest version is Python 3.8. Then scroll to the bottom and select *Windows x86-64 executable installer*.

Note

If your system has a 32-bit processor, then you should choose the 32-bit installer. If you aren't sure if your computer is 32-bit or 64-bit, stick with the 64-bit installer mentioned above.

Step 2: Run the Installer

Run the installer by double-clicking on the downloaded file. You should see the following window:



Important

Make sure you check the box that says *Add Python 3.x to PATH* as shown to ensure that the install places the interpreter in your execution path.

If you install Python without checking this box, you can run the installer again and select it.

Click **Install Now** to install Python 3. Wait for the installation to finish, and then continue to open IDLE.

Open IDLE

You can open IDLE in two steps:

1. Click on the start menu and locate the *Python 3.8* folder.
2. Open the folder and select *IDLE (Python 3.8)*.

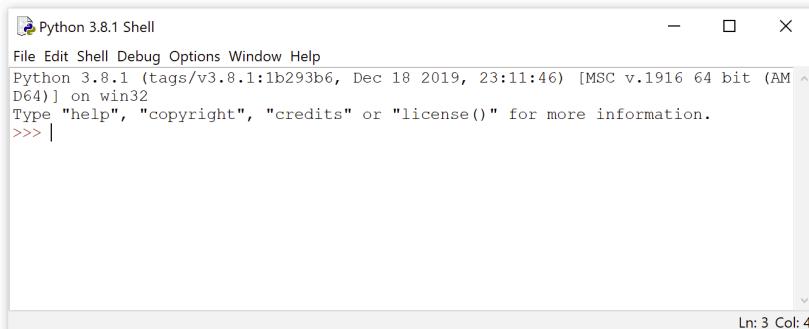
Note

We recommend using IDLE to follow along with this book.

You may use a different code editor if you prefer. However, some chapters, such as Chapter 7: *Finding And Fixing Code Bugs*, contain material specific to IDLE.

IDLE opens a **Python shell** in a new window. The Python shell is an interactive environment that allows you to type in some Python code and execute it immediately. It is a great way to get started with Python!

The Python shell window looks like this:



At the top of the window, you can see the version of Python that is running and some information about the operating system. If you see a version less than 3.7, you may need to revisit the installation instructions in the previous section.

The `>>>` symbol that you see is called a **prompt**. Whenever you see this, it means that Python is waiting for you to give it some instructions.

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-2

Now that you have Python installed, let's get straight into writing your first Python program! Go ahead and move on to Chapter 3.

[Leave feedback on this section »](#)

2.3 macOS

Follow these steps to install Python 3 and open IDLE on macOS.

Important

The code in this book is only tested against Python installed as described in this section.

Be aware that if you have installed Python through some other means, such as Homebrew or Anaconda Python, you may encounter problems when running the code examples.

Install Python

Most macOS machines come with Python 2 installed. You'll want to install the latest version of Python 3. You can do this by downloading an installer from the python.org website.

Step 1: Download the Python 3 Installer

Open a browser window and navigate to the download page for macOS at python.org.

Underneath the heading at the top that says *Python Releases for macOS*, click on the link for the *Latest Python 3 Release - Python 3.x.x*. As of this writing, the latest version is Python 3.8. Then scroll to the bottom of the page and select *macOS 64-bit/32-bit installer*. This starts the download.

Step 2: Run the Installer

Run the installer by double-clicking on the downloaded file. You should see the following window:



1. Press the **Continue** button a few times until you are asked to agree to the software license agreement. Then click **Agree**. You are shown a window that tells you where Python will be installed and how much space it will take.
2. You most likely don't want to change the default location, so go ahead and click **Install** to start the installation. The Python installer will tell you when it is finished copying files.

3. Click **Close** to close the installer window. Now that Python is installed, you can open up IDLE and get ready to write your first Python program.

Open IDLE

You can open IDLE in three steps:

1. Open Finder and click on *Applications*.
2. Locate the *Python 3.8* folder and double-click on it.
3. Double-click on the IDLE icon.

You may also open IDLE using the Spotlight search feature. Press **Cmd** + **Spacebar** to open the Spotlight search, type the word `idle`, then press **Return** to open IDLE.

Note

We recommend using IDLE to follow along with this book.

You may use a different code editor if you prefer. However, some chapters, such as Chapter 7: *Finding And Fixing Code Bugs*, contain material specific to IDLE.

IDLE opens a **Python shell** in a new window. The Python shell is an interactive environment that allows you to type in some Python code and execute it immediately. It is a great way to get started with Python!

The Python shell window looks like this:

```
Python 3.8.1 (v3.8.1:1b293b6006, Dec 18 2019, 14:08:53)
[Clang 6.0 (clang-600.0.57)] on darwin
Type "help", "copyright", "credits" or "license()" for more information.
>>> |
```

Ln: 4 Col: 4

At the top of the window, you can see the version of Python that is running and some information about the operating system. If you see a version less than 3.7, you may need to revisit the installation instructions in the previous section.

The `>>>` symbol that you see is called a **prompt**. Whenever you see this, it means that Python is waiting for you to give it some instructions.

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-2

Now that you have Python installed, let's get straight into writing your first Python program! Go ahead and move on to Chapter 3.

[Leave feedback on this section »](#)

2.4 Ubuntu Linux

Follow these steps to install Python 3 and open IDLE on Ubuntu Linux.

Important

The code in this book is only tested against Python installed as described in this section.

Be aware that if you have installed Python through some other means, such as Anaconda Python, you may encounter problems when running the code examples.

Install Python

There is a good chance your Ubuntu distribution has Python installed already, but it probably won't be the latest version, and it may be Python 2 instead of Python 3.

To find out what version(s) you have, open a terminal window and try the following commands:

```
$ python --version  
$ python3 --version
```

One or more of these commands should respond with a version, as below (your version number may vary):

```
$ python3 --version  
Python 3.8.1
```

If the version shown is Python 2.x or a version of Python 3 that is less than 3.8, then you want to install the latest version. How you install Python on Ubuntu depends on which version of Ubuntu you are running. You can determine your local Ubuntu version by running the following command:

```
$ lsb_release -a
No LSB modules are available.
Distributor ID: Ubuntu
Description:    Ubuntu 18.04.1 LTS
Release:        18.04
Codename:       bionic
```

Look at the version number next to `Release` in the console output, and follow the corresponding instructions below.

Ubuntu 18.04+

Ubuntu version 18.04 does not come with Python 3.8 by default, but it is in the Universe repository. You can install it with the following commands in the *Terminal* application:

```
$ sudo apt-get update
$ sudo apt-get install python3.8 idle-python3.8
```

Ubuntu 17 and lower

For Ubuntu versions 17 and lower, Python 3.8 is not in the Universe repository. You need to get it from a Personal Package Archive (PPA). To install Python from the “[deadsnakes](#)” PPA, run the following commands in the *Terminal* application:

```
$ sudo add-apt-repository ppa:deadsnakes/ppa
$ sudo apt-get update
$ sudo apt-get install python3.8 idle-python3.8
```

You can check that the correct version of Python was installed by running `python3 --version`. If you see a version number less than 3.7, you may need to type `python3.8 --version`. Now you are ready to open IDLE and get ready to write your first Python program.

Open IDLE

You can open IDLE from the command line by typing the following:

```
$ idle-python3.8
```

On some Linux installations, you can open IDLE with the following shortened command:

```
$ idle3
```

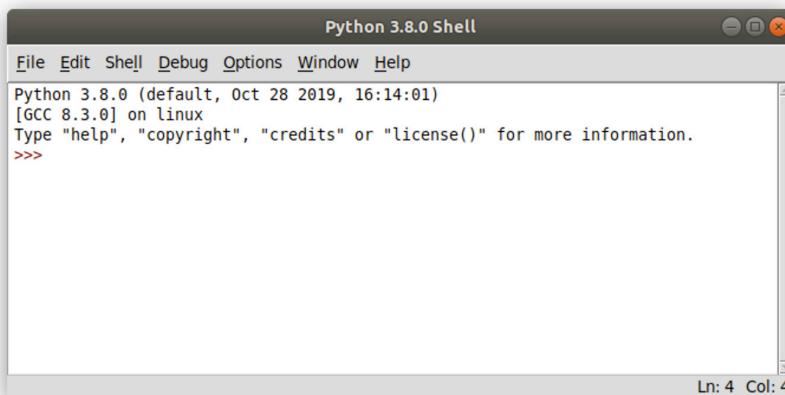
Note

We recommend using IDLE to follow along with this book.

You may use a different code editor if you prefer. However, some chapters, such as Chapter 7: *Finding And Fixing Code Bugs*, contain material specific to IDLE.

IDLE opens a **Python shell** in a new window. The Python shell is an interactive environment that allows you to type in some Python code and execute it immediately. It is a great way to get started with Python!

The Python shell window looks like this:



At the top of the window, you can see the version of Python that is

running and some information about the operating system. If you see a version less than 3.8, you may need to revisit the installation instructions in the previous section.

Important

If you opened IDLE with the `idle3` command and see a version less than 3.7 displayed in the Python shell window, then you will need to open IDLE with the `idle-python3.8` command.

The `>>>` symbol that you see in the IDLE window is called a **prompt**. Whenever you see this, it means that Python is waiting for you to give it some instructions.

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-2

Now that you have Python installed, let's get straight into writing your first Python program! Go ahead and move on to Chapter 3.

[Leave feedback on this section »](#)

Chapter 3

Your First Python Program

Now that you have the latest version of Python installed on your computer, it's time to start coding!

In this chapter, you will:

- Write your first Python script
- Learn what happens when you run a script with an error
- Learn how to declare a variable and inspect its value
- Learn how to write comments

Ready to begin your Python journey? Let's go!

[Leave feedback on this section »](#)

3.1 Write a Python Script

If you don't have IDLE open already, go ahead and open it. There are two main windows that you will work with in IDLE: the **interactive window**, which is the one that opens when you start IDLE, and the **script window**.

You can type code into both the interactive and script windows. The difference between the two is how the code is executed. In this section,

you will write your first Python program and learn how to run it in both windows.

The Interactive Window

The interactive window contains a **Python shell**, which is a textual user interface used to interact with the Python language. Hence the name “interactive window.”

When you first open IDLE, the text displayed looks something like this:

```
Python 3.8.1 (tags/v3.8.1:1b293b6, Dec 18 2019, 22:39:24)
[MSC v.1916 32 bit (Intel)] on win32
Type "help", "copyright", "credits" or "license" for more information.
>>>
```

The first line tells you what version of Python is running. In this case, IDLE is running Python 3.8.1. The second and third lines give some information about the operating system and some commands you can use to get more information about Python.

The `>>>` symbol in the last line is called the **prompt**. This is where you will type in your code. Go ahead and type `1 + 1` at the prompt and press Enter.

When you hit `Enter`, Python evaluates the expression, displays the result 2, and then prompts you for more input:

```
>>> 1 + 1
2
>>>
```

Notice that the Python prompt `>>>` appears again after your result. Python is ready for more instructions! Every time you run some code, a new prompt appears directly below the output.

The sequence of events in the interactive window can be described as a loop with three steps:

1. First, Python reads the code entered at the prompt.
2. Then the code is evaluated.
3. Finally, the output is printed in the window and a new prompt is displayed.

This loop is commonly referred to as a **Read-Evaluate-Print Loop**, or **REPL**. Python programmers sometimes refer the Python shell as a “Python REPL”, or just “the REPL” for short.

Note

From this point on, the final `>>>` prompt displayed after executing code in the interactive window is excluded from code examples.

Let's try something a little more interesting than adding two numbers. A rite of passage for every programmer is writing their first “Hello, world” program that prints the phrase “Hello, world” on the screen.

To print text to the screen in Python, you use the `print()` function. A **function** is a bit of code that typically takes some input, called an **argument**, does something with that input, and produces some output, called the **return value**.

Loosely speaking, functions in code work like mathematical functions. For example, the mathematical function $A(r)=\pi r^2$ takes the radius r of a circle as input and produces the area of the circle as output.

Important

The analogy to mathematical functions has some problems, though, because code functions can have **side effects**. A side effect occurs anytime a function performs some operation that changes something about the program or the computer running the program.

For example, you can write a function in Python that takes someone's name as input, stores the name in a file on the computer, and then outputs the path to the file with the name in it. The operation of saving the name to a file is a side effect of the function.

You'll learn more about functions, including how to write your own, in Chapter 6.

Python's `print()` function takes some text as input and then displays that text on the screen. To use `print()`, type the word `print` at the prompt in the interactive window, followed by the text "Hello, world" inside of parentheses:

```
>>> print("Hello, world")
Hello, world
```

Here "Hello, world" is the argument that is being **passed** to `print()`. "Hello, world" must be written with quotation marks so that Python interprets it as text and not something else.

Note

As you type code into the interactive window, you may notice that the font color changes for certain parts of the code. IDLE **highlights** parts of your code in different colors to help make it easier for you to identify what the different parts are.

By default, built-in functions, such as `print()` are displayed in purple, and text is displayed in green.

The interactive window can execute only a single line of code at a time. This is useful for trying out small code examples and exploring the Python language, but it has a major limitation. Code must be entered in by a person one line at a time!

Alternatively, you can store some Python code in a text file and then execute all of the code in the file with a single command. The code in the file is called a **script**, and files containing Python scripts are called **script files**.

Script files are nice not only because they make it easier to run a program, but also because they can be shared with other people so that they can run your program, too.

The Script Window

Scripts are written using IDLE's script window. You can open the script window by selecting **File > New File** from the menu at the top of the interactive window.

Notice that when the script window opens, the interactive window stays open. Any output generated by code run in the script window is displayed in the interactive window, so you may want to rearrange the two windows so that you can see both of them at the same time.

In the script window, type in the same code you used to print "Hello, world" in the interactive window:

```
print("Hello, world")
```

Just like the interactive window, code typed into the script window is highlighted.

Important

When you write code in a script, you do not need to include the `>>>` prompt that you see in IDLE's interactive window. Keep this in mind if you copy and paste code from examples that show the REPL prompt.

Remember, though, that it's not recommended that you copy and paste examples from the book. Typing each example in yourself really pays off!

Before you can run your script, you must save it. From the menu at the top of the window, select `File > Save` and save the script as `hello_world.py`. The `.py` file extension is the conventional extension used to indicate that a file contains Python code.

In fact, if you save your script with any extension other than `.py`, the code highlighting will disappear and all the text in the file will be displayed in black. IDLE will only highlight Python code when it is stored in a `.py` file.

Once the script is saved, all you have to do to run the program is select `Run > Run Module` from the script window and you'll see `Hello, world` appear in the interactive window:

```
Hello, world
```

Note

You can also press `F5` to run a script from the script window.

Every time you run a script you will see something like the following output in the interactive window:

```
>>> ===== RESTART =====
```

This is IDLE's way of separating output from distinct runs of a script. Otherwise, if you run one script after another, it may not be clear what

output belongs to which script.

To open an existing script in IDLE, select `File > Open...` from the menu in either the script window or the interactive window. Then browse for and select the script file you want to open. IDLE opens scripts in a new script window, so you can have several scripts open at a time.

Note

Double-clicking on a `.py` file from a file manager, such as Windows Explorer, does execute the script in a new window. However, the window is closed immediately when the script is done running—often before you can even see what happened.

To open the file in IDLE so that you can run it and see the output, you can right-click on the file icon (`Ctrl`+`Click` on macOS) and choose to `Edit with IDLE`.

[Leave feedback on this section »](#)

3.2 Mess Things Up

Everybody makes mistakes—especially while programming! In case you haven’t made any mistakes yet, let’s get a head start on that and mess something up on purpose to see what happens.

Mistakes made in a program are called **errors**, and there are two main types of errors you’ll experience:

1. Syntax errors
2. Run-time errors

In this section you’ll see some examples of code errors and learn how to use the output Python displays when an error occurs to understand what error occurred and which piece of code caused it.

Syntax Errors

In loose terms, a **syntax error** occurs when you write some code that isn't allowed in the Python language. You can create a syntax error by changing the contents of the `hello_world.py` script from the last section to the following:

```
print("Hello, world)
```

In this example, the double quotation mark at the end of "Hello, world" has been removed. Python won't be able to tell where the string of text ends. Save the altered script and then try to run it. What happens?

The code won't run! IDLE displays an alert box with the following message:

```
EOL while scanning string literal.
```

EOL stands for **End Of Line**, so this message tells you that Python read all the way to the end of the line without finding the end of something called a string literal.

A **string literal** is text contained in-between two double quotation marks. The text "Hello, world" is an example of a string literal.

Note

For brevity, string literals are often referred to as **strings**, although the term "string" technically has a more general meaning in Python. You will learn more about strings in Chapter 4.

Back in the script window, notice that the line containing with "Hello, world is highlighted in red. This handy features helps you quickly find which line of code caused the syntax error.

Run-time Errors

IDLE catches syntax errors before a program starts running, but some errors can't be caught until a program is executed. These errors are

known as **run-time errors** because they only occur at the time that a program is run.

To generate a run-time error, change the code in `hello_world.py` to the following:

```
print(Hello, world)
```

Now both quotation marks from the phrase "Hello, world" have been removed. Did you notice how the text color changes to black when you removed the quotation marks? IDLE no longer recognizes `Hello, world` as a string.

What do you think happens when you run the script? Try it out and see!

Some red text is displayed in the interactive window:

```
Traceback (most recent call last):
  File "/home/hello_world.py", line 1, in <module>
    print(Hello, world)
NameError: name 'Hello' is not defined
```

What happened? While trying to execute the program Python **raised** an error. Whenever an error occurs, Python stops executing the program and displays the error in IDLE's interactive window.

The text that gets displayed for an error is called a **traceback**. Tracebacks give you some useful information about the error. The traceback above tells us all of the following:

- The error happened on line 1 of the `hello_world.py`.
- The line that generated the error was: `print(Hello, world)`.
- A `NameError` occurred.
- The specific error was `name 'Hello' is not defined`

The quotation marks around `Hello, world` are missing, so Python doesn't understand that it is a string of text. Instead, Python thinks

that `Hello` and `world` are the names of something else in the code. Since names `Hello` and `world` haven't been defined anywhere, the program crashes.

In the next section, you'll see how to define names for values in your code. Before you move on though, you can get some practice with syntax errors and run-time errors by working on the review exercises.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Write a script that IDLE won't let you run because it has a syntax error.
2. Write a script that only crashes your program once it is already running because it has a run-time error.

[Leave feedback on this section »](#)

3.3 Create a Variable

In Python, **variables** are names that can be assigned a value and used to reference that value throughout your code. Variables are fundamental to programming for two reasons:

1. **Variables keep values accessible:** For example, the result of some time-consuming operation can be assigned to a variable so that the operation does not need to be performed each time you need to use the result.
2. **Variables give values context:** The number 28 could mean lots of different things, such as the number of students in a class, or the number of times a user has accessed a website, and so on. Naming the value 28 something like `num_students` makes the meaning of the value clear.

In this section, you'll learn how to use variables in your code, as well as some of the conventions Python programmers follow when choosing names for variables.

The Assignment Operator

Values are assigned to a variable using a special symbol = called the **assignment operator**. An **operator** is a symbol, like = or +, that performs some operation on one or more values.

For example, the + operator takes two numbers, one to the left of the operator and one to the right, and adds them together. Likewise, the = operator takes a value to the right of the operator and assigns it to the name on the left of the operator.

To see the assignment operator in action, let's modify the "Hello, world" program you saw in the last section. This time, we'll use a variable to store some text before printing it to the screen:

```
>>> phrase = "Hello, world"  
>>> print(phrase)  
Hello, world
```

In the first line, a variable named `phrase` is created and assigned the value "Hello, world" using the = operator. The string "Hello, world" that was originally used inside of the parentheses in the `print()` function is replaced with the variable `phrase`.

The output `Hello, world` is displayed when you execute `print(phrase)` because Python looks up the name `phrase` and finds it has been assigned the value "Hello, world".

If you hadn't executed `phrase = "Hello, world"` before executing `print(phrase)`, you would have seen a `NameError` like you did when trying to execute `print(Hello, world)` in the previous section.

Note

Although = looks like the equals sign from mathematics, it has a different meaning in Python. Distinguishing the = operator from the equals sign is important, and can be a source of frustration for beginner programmers.

Just remember, whenever you see the = operator, whatever is to the right of it is being assigned to a variable on the left.

Variable names are **case-sensitive**, so a variable named `phrase` is distinct from a variable named `Phrase` (note the capital P). For instance, the following code produces a `NameError`:

```
>>> phrase = "Hello, world"
>>> print(Phrase)
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
NameError: name 'Phrase' is not defined
```

When you run into trouble with the code examples in this book, be sure to double-check that every character in your code—including spaces—exactly matches the examples. Computers can't use common sense to interpret what you meant to say, so being *almost* correct won't get a computer to do the right thing!

Rules for Valid Variable Names

Variable names can be as long or as short as you like, but there are a couple of rules that you must follow. Variable names can only contain uppercase and lowercase letters (A–Z, a–z), digits (0–9), and underscores (_). However, variable names cannot begin with a digit.

For example, `phrase`, `string1`, `_a1p4a`, and `list_of_names` are all valid variable names, but `9lives` is not.

Note

Python variable names can contain many different valid Unicode characters. **Unicode** is a standard for digitally representing text used in most of the world's writing systems.

That means variable names can contain letters from non-English alphabets, such as decorated letters like é and ü, and even Chinese, Japanese, and Arabic symbols.

However, not every system can display decorated characters, so it is a good idea to avoid them if your code is going to be shared with people in many different regions.

You can learn more about Unicode on [Wikipedia](#). Python's support for Unicode is covered in the [official Python documentation](#).

Just because a variable name is valid doesn't necessarily mean that it is a good name. Choosing a good name for a variable can be surprisingly difficult. However, there are some guidelines that you can follow to help you choose better names.

Descriptive Names Are Better Than Short Names

Descriptive variable names are essential, especially for complex programs. Often, descriptive names require using multiple words. Don't be afraid to use long variable names.

In the following example, the value 3600 is assigned to the variable s:

```
s = 3600
```

The name s is totally ambiguous. Using a full word makes it a lot easier to understand what the code means:

```
seconds = 3600
```

seconds is a better name than s because it provides more context. But

it still doesn't convey the full meaning of the code. Is 3600 the number of seconds it takes for some process to finish, or the length of a movie? There's no way to tell.

The following name leaves no doubt about what the code means:

```
seconds_per_hour = 3600
```

When you read the above code, there is no question that 3600 is the number of seconds in one hour. Although `seconds_per_hour` takes longer to type than both the single letter `s` and the word `seconds`, the pay-off in clarity is massive.

Although naming variables descriptively means using longer variable names, you should avoid names that are excessively long. What “excessively long” really means is subjective, but a good rule of thumb is to keep variable names to fewer than three or four words.

Python Variable Naming Conventions

In many programming languages, it is common to write variable names in **camelCase** like `numStudents` and `listOfNames`. The first letter of every word, except the first, is capitalized, and all other letters are lowercase. The juxtaposition of lower-case and upper-case letters look like humps on a camel.

In Python, however, it is more common to write variable names in **snake case** like `num_students` and `list_of_names`. Every letter is lowercase, and each word is separated by an underscore.

While there is no hard-and-fast rule mandating that you write your variable names in snake case, the practice is codified in a document called **PEP 8**, which is widely regarded as the official style guide for writing Python.

Following the standards outlined in PEP 8 ensures that your Python code is readable by a large number of Python programmers. This makes sharing and collaborating on code easier for everyone involved.

Note

All of the code examples in this course follow PEP 8 guidelines, so you will get a lot of exposure to what Python code that follows standard formatting guidelines looks like.

In this section you learned how to create a variable, rules for valid variable names, and some guidelines for choosing good variable names. Next, you will learn how to inspect a variable's value in IDLE's interactive window.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Using the interactive window, display some text on the screen by using the `print()` function.
2. Using the interactive window, display a string of text by saving the string to a variable, then reference the string in a `print()` function using the variable name.
3. Do each of the first two exercises again by first saving your code in a script and running it.

[Leave feedback on this section »](#)

3.4 Inspect Values in the Interactive Window

You have already seen how to use `print()` to display a string that has been assigned to a variable. There is another way to display the value of a variable when you are working in the Python shell.

Type the following into IDLE's interactive window:

```
>>> phrase = "Hello, world"  
>>> phrase
```

When you press Enter after typing `phrase` a second time, the following output is displayed:

```
'Hello, world'
```

Python prints the string "Hello, world", and you didn't have to type `print(phrase)!`

Now type the following:

```
>>> print(phrase)
```

This time, when you hit Enter you see:

```
Hello, world
```

Do you see the difference between this output and the output of simply typing `phrase`? It doesn't have any single quotes surrounding it. What's going on here?

When you type `phrase` and press Enter, you are telling Python to **inspect** the variable `phrase`. The output displayed is a useful representation of the value assigned to the variable.

In this case, `phrase` is assigned the string "Hello, world", so the output is surrounded with single quotes to indicate that `phrase` is a string.

On the other hand, when you `print()` a variable, Python displays a more human-readable representation of the variable's value. For strings, both ways of being displayed are human-readable, but this is not the case for every type of value.

Sometimes, both printing and inspecting a variable produces the same output:

```
>>> x = 2
>>> x
2
>>> print(x)
2
```

Here, `x` is assigned to the number 2. Both the output of `print(x)` and inspecting `x` is not surrounded with quotes, because 2 is a number and not a string.

Inspecting a variable, instead of printing it, is useful for a couple of reasons. You can use it to display the value of a variable without typing `print()`. More importantly, though, inspecting a variable usually gives you more useful information than `print()` does.

Suppose you have two variables: `x = 2` and `y = "2"`. In this case, `print(x)` and `print(y)` both display the same thing. However, inspecting `x` and `y` shows the difference between the each variable's value:

```
>>> x = 2
>>> y = "2"
>>> print(x)
2
>>> print(y)
2
>>> x
2
>>> y
'2'
```

The key takeaway here is that `print()` displays a readable representation of a variable's value, while inspection provides additional information about the type of the value.

You can inspect more than just variables in the Python shell. Check out what happens when you type `print` and hit Enter:

```
>>> print  
<built-in function print>
```

Keep in mind that you can only inspect variables in a Python shell. For example, save and run the following script:

```
phrase = "Hello, world"  
phrase
```

The script executes without any errors, but no output is displayed! Throughout this book, you will see examples that use the interactive window to inspect variables.

[Leave feedback on this section »](#)

3.5 Leave Yourself Helpful Notes

Programmers often read code they wrote several months ago and wonder “What the heck does this do?” Even with descriptive variable names, it can be difficult to remember why you wrote something the way you did when you haven’t looked at it for a long time.

To help avoid this problem, you can leave comments in your code. **Comments** are lines of text that don’t affect the way the script runs. They help to document what’s supposed to be happening.

In this section, you will learn three ways to leave comments in your code. You will also learn some conventions for formatting comments, as well as some pet peeves regarding their over-use.

How to Write a Comment

The most common way to write a comment is to begin a new line in your code with the `#` character. When your code is run, any lines starting with `#` are ignored. Comments that start on a new line are called **block comments**.

You can also write **in-line comments**, which are comments that appear on the same line as some code. Just put a # at the end of the line of code, followed by the text in your comment.

Here is an example of the `hello_world.py` script with both kinds of comments added in:

```
# This is a block comment.

phrase = "Hello, world."
print(phrase) # This is an in-line comment.
```

The first line doesn't do anything, because it starts with a #. Likewise, `print(phrase)` is executed on the last line, but everything after the # is ignored.

Of course, you can still use the # symbol inside of a string. For instance, Python won't mistake the following for the start of a comment:

```
print("#1")
```

In general, it's a good idea to keep comments as short as possible, but sometimes you need to write more than will reasonably fit on a single line. In that case, you can continue your comment on a new line that also begins with a # symbol:

```
# This is my first script.
# It prints the phrase "Hello, world."
# The comments are longer than the script!

phrase = "Hello, world."
print(phrase)
```

Besides leaving yourself notes, comments can also be used to **comment out** code while you're testing a program. In other words, adding a # at the beginning of a line of code lets you run your program as if that line of code didn't exist without having to delete any code.

To comment out a section of code in IDLE highlight one or more lines

to be commented and press:

- **Windows:** `Alt + 3`
- **macOS:** `Ctrl + 3`
- **Ubuntu Linux:** `Ctrl + D`

Two # symbols are inserted at the beginning of each line. This doesn't follow PEP 8 comment formatting conventions, but it gets the job done!

To un-comment out your code and remove the # symbols from the beginning of each line, highlight the code that is commented out and press:

- **Windows:** `Alt + 4`
- **macOS:** `Ctrl + 4`
- **Ubuntu Linux:** `Ctrl + Shift + D`

Now let's look at some common conventions regarded code comments.

Conventions and Pet Peeves

According to [PEP 8](#), comments should always be written in complete sentences with a single space between the # and the first word of the comment:

```
# This comment is formatted to PEP 8.  
  
#don't do this
```

For in-line comments, PEP 8 recommends at least two spaces between the code and the # symbol:

```
phrase = "Hello, world" # This comment is PEP 8 compliant.  
print(phrase) # This comment isn't.
```

A major pet peeve among programmers are comments that describe

what is already obvious from reading the code. For example, the following comment is unnecessary:

```
# Print "Hello, world"  
print("Hello, world")
```

No comment is needed in this example because the code itself explicitly describes what is being done. Comments are best used to clarify code that may not be easy to understand, or to explain why something is done a certain way.

In general, PEP 8 recommends that comments be used sparingly. Use comments only when they add value to your code by making it easier to understand *why* something is done a certain way. Comments that describe *what* something does can often be avoided by using more descriptive variable names.

[Leave feedback on this section »](#)

3.6 Summary and Additional Resources

In this chapter, you wrote and executed your first Python program! You wrote a small program that displays the text "Hello, world" using the `print()` function.

You were introduced to three concepts:

1. **Variables** give names to values in your code using the assignment operator (`=`)
2. **Errors**, such as syntax errors and run-time errors, are raised whenever Python can't execute your code. They are displayed in IDLE's interactive window in the form of a traceback.
3. **Comments** are lines of code that don't get executed and serve as documentation for yourself and other programmers that need to read your code.

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-3

Additional Resources

To learn more, check out the following resources:

- [11 Beginner Tips for Learning Python Programming](#)
- [Writing Comments in Python \(Guide\)](#)
- Recommended resources on realpython.com

[Leave feedback on this section »](#)

Chapter 4

Strings and String Methods

Many programmers, regardless of their specialty, deal with text on a daily basis. For example, web developers work with text that gets input from web forms. Data scientists process text to extract data and perform things like sentiment analysis, which can help identify and classify opinions in a body of text.

Collections of text in Python are called **strings**. Special functions called **string methods** are used to manipulate strings. There are string methods for changing a string from lowercase to uppercase, removing whitespace from the beginning or end of a string, or replacing parts of a string with different text, and many more.

In this chapter, you will learn how to:

- Manipulate strings with string methods
- Work with user input
- Deal with strings of numbers
- Format strings for printing

Let's get started!

[Leave feedback on this section »](#)

4.1 What is a String?

In Chapter 3, you created the string "Hello, world" and printed it in IDLE's interactive window using the `print()` function. In this section, you'll get a deeper look into what exactly a string is and the various ways you can create them in Python.

The String Data Type

Strings are one of the fundamental Python data types. The term **data type** refers to what kind of data a value represents. Strings are used to represent text.

Note

There are several other data types built-in to Python. For example, you'll learn about numerical data types in Chapter 5, and Boolean types in Chapter 8.

We say that strings are a **fundamental** data type because they can't be broken down into smaller values of a different type. Not all data types are fundamental. You'll learn about compound data types, also known as **data structures**, in Chapter 9.

The string data type has a special abbreviated name in Python: `str`. You can see this by using the `type()` function, which is used to determine the data type of a given value.

Type the following into IDLE's interactive window:

```
>>> type("Hello, world")
<class 'str'>
```

The output `<class 'str'>` indicates that the value "Hello, world" is an instance of the `str` data type. That is, "Hello, world" is a string.

Note

For now, you can think of the word “class” as a synonym for “data type,” although it actually refers to something more specific. You’ll see just what a class is in Chapter 10.

`type()` also works for values that have been assigned to a variable:

```
>>> phrase = "Hello, world"  
>>> type(phrase)  
<class 'str'>
```

Strings have three properties that you’ll explore in the coming sections:

1. Strings contain **characters**, which are individual letters or symbols.
2. Strings have a **length**, which is the number of characters contained in the string.
3. Characters in a string appear in a **sequence**, meaning each character has a numbered position in the string.

Let’s take a closer look at how strings are created.

String Literals

As you’ve already seen, you can create a string by surrounding some text with quotation marks:

```
string1 = 'Hello, world'  
string2 = "1234"
```

Either single quotes (`string1`) or double quotes (`string2`) can be used to create a string, as long as both quotation marks are the same type.

Whenever you create a string by surrounding text with quotation marks, the string is called a **string literal**. The name indicates that the string is literally written out in your code. All of the strings you

have seen thus far are string literals.

Note

Not every string is a string literal. For example, a string captured as user input isn't a string literal because it isn't explicitly written out in the program's code.

You'll learn how to work with user input in section 4 of this chapter.

The quotes surrounding a string are called **delimiters** because they tell Python where a string begins and where it ends. When one type of quotes is used as the delimiter, the other type of quote can be used inside of the string:

```
string3 = "We're #1!"  
string4 = 'I said, "Put it over by the llama.''
```

After Python reads the first delimiter, all of the characters after it are considered a part of the string until a second matching delimiter is read. This is why you can use a single quote in a string delimited by double quotes and vice versa.

If you try to use double quotes inside of a string that is delimited by double quotes, you will get an error:

```
>>> text = "She said, "What time is it?""  
      File "<stdin>", line 1  
      text = "She said, "What time is it?""  
                           ^  
SyntaxError: invalid syntax
```

Python throws a `SyntaxError` because it thinks that the string ends after the second `"` and doesn't know how to interpret the rest of the line.

Note

A common pet peeve among programmers is the use of mixed quotes as delimiters. When you work on a project, it's a good idea to use only single quotes or only double quotes to delimit every string.

Keep in mind that there isn't really a right or wrong choice! The goal is to be consistent, because consistency helps make your code easier to read and understand.

Strings can contain any valid Unicode character. For example, the string "We're #1!" contains the pound sign (#) and "1234" contains numbers. "xPýthønx" is also a valid Python string!

Determine the Length of a String

The number of characters contained in a string, including spaces, is called the **length** of the string. For example, the string "abc" has a length of 3, and the string "Don't Panic" has a length of 11.

To determine a string's length, you use Python's built-in `len()` function. To see how it works, type the following into IDLE's interactive window:

```
>>> len("abc")
3
```

You can also use `len()` to get the length of a string that's assigned to a variable:

```
>>> letters = "abc"
>>> num_letters = len(letters)
>>> num_letters
3
```

First, the string "abc" is assigned to the variable `letters`. Then `len()` is used to get the length of `letters` and this value is assigned to the `num_letters` variable. Finally, the value of `num_letters`, which is 3, is

displayed.

Multiline Strings

The [PEP 8](#) style guide recommends that each line of Python code contain no more than 79 characters—including spaces.

Note

PEP 8's 79-character line-length is recommended because, among other things, it makes it easier to read two files side-by-side. However, many Python programmers believe forcing each line to be at most 79 characters sometimes makes code harder to read.

In this book we will strictly follow PEP 8's recommended line-length. Just know that you will encounter lots of code in the real world with longer lines.

Whether you decide to follow PEP 8, or choose a larger number of characters for your line-length, you will sometimes need to create string literals with more characters than your chosen limit.

To deal with long strings, you can break the string up across multiple lines into a **multiline string**. For example, suppose you need to fit the following text into a string literal:

“This planet has—or rather had—a problem, which was this: most of the people living on it were unhappy for pretty much of the time. Many solutions were suggested for this problem, but most of these were largely concerned with the movements of small green pieces of paper, which is odd because on the whole it wasn’t the small green pieces of paper that were unhappy.”

— Douglas Adams, *The Hitchhiker’s Guide to the Galaxy*

This paragraph contains far more than 79 characters, so any line of code containing the paragraph as a string literal violates PEP 8. So, what do you do?

There are a couple of ways to tackle this. One way is to break the string up across multiple lines and put a backslash (\) at the end of all but the last line. To be PEP 8 compliant, the total length of the line, including the backslash, must be 79 characters or less.

Here's how you could write the paragraph as a multiline string using the backslash method:

```
paragraph = "This planet has - or rather had - a problem, which was \
this: most of the people living on it were unhappy for pretty much \
of the time. Many solutions were suggested for this problem, but \
most of these were largely concerned with the movements of small \
green pieces of paper, which is odd because on the whole it wasn't \
the small green pieces of paper that were unhappy."
```

Notice that you don't have to close each line with a quotation mark. Normally, Python would get to the end of the first line and complain that you didn't close the string with a matching double quote. With a backslash at the end, however, you can keep writing the same string on the next line.

When you `print()` a multiline string that is broken up by backslashes, the output displayed on a single line:

```
>>> long_string = "This multiline string is \
displayed on one line"
>>> print(long_string)
This multiline string is displayed on one line
```

Multiline strings can also be created using triple quotes as delimiters (""" or '''). Here is how you might write a long paragraph using this approach:

```
paragraph = """This planet has - or rather had - a problem, which was  
this: most of the people living on it were unhappy for pretty much  
of the time. Many solutions were suggested for this problem, but  
most of these were largely concerned with the movements of small  
green pieces of paper, which is odd because on the whole it wasn't  
the small green pieces of paper that were unhappy."""
```

Triple-quoted strings preserve whitespace. This means that running `print(paragraph)` displays the string on multiple lines just like it is in the string literal, including newlines. This may or may not be what you want, so you'll need to think about the desired output before you choose how to write a multiline string.

To see how whitespace is preserved in a triple-quoted string, type the following into IDLE's interactive window:

```
>>> print("""An example of a  
...     string that spans across multiple lines  
...             that also preserves whitespace.""")  
An example of a  
    string that spans across multiple lines  
        that also preserves whitespace.
```

Notice how the second and third lines in the output are indented exactly the same way they are in the string literal.

Note

Triple-quoted strings have a special purpose in Python. They are used to document code. You'll often find them at the top of a .py with a description of the code's purpose. They are also used to document custom functions.

When used to document code, triple-quoted strings are called **docstrings**. You'll learn more about docstrings in Chapter 6.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Print a string that uses double quotation marks inside the string.
2. Print a string that uses an apostrophe inside the string.
3. Print a string that spans multiple lines, with whitespace preserved.
4. Print a string that is coded on multiple lines but displays on a single line.

[Leave feedback on this section »](#)

4.2 Concatenation, Indexing, and Slicing

Now that you know what a string is and how to declare string literals in your code, let's explore some of the things you can do with strings.

In this section, you'll learn about three basic string operations:

1. Concatenation, which joins two strings together
2. Indexing, which gets a single character from a string
3. Slicing, which gets several characters from a string at once

Let's dive in!

String Concatenation

Two strings can be combined, or **concatenated**, using the `+` operator:

```
>>> string1 = "abra"
>>> string2 = "cadabra"
>>> magic_string = string1 + string2
>>> magic_string
```

```
'abracadabra'
```

In this example, string concatenation occurs on the third line. `string1` and `string2` are concatenated using `+` and the result is assigned to the variable `magic_string`. Notice that the two strings are joined without any whitespace between them.

You can use string concatenation to join two related strings, such as joining a first and last name into a full name:

```
>>> first_name = "Arthur"
>>> last_name = "Dent"
>>> full_name = first_name + " " + last_name
>>> full_name
'Arthur Dent'
```

Here string concatenation occurs twice on the same line. `first_name` is concatenated with `" "`, resulting in the string `"Arthur "`. Then this result is concatenated with `last_name` to produce the full name `"Arthur Dent"`.

String Indexing

Each character in a string has a numbered position called an **index**. You can access the character at the *Nth* position by putting the number *N* in between two square brackets (`[` and `]`) immediately after the string:

```
>>> flavor = "apple pie"
>>> flavor[1]
'p'
```

`flavor[1]` returns the character at position 1 in `"apple pie"`, which is `p`. Wait, isn't `a` the first character of `"apple pie"`?

In Python—and most other programming languages—counting always starts at zero. To get the character at the beginning of a string, you need to access the character at position 0:

```
>>> flavor[0]  
'a'
```

Note

Forgetting that counting starts with zero and trying to access the first character in a string with the index 1 results in an **off-by-one error**.

Off-by-one errors are a common source of frustration for both beginning and experienced programmers alike!

The following figure shows the index for each character of the string "apple pie":

	a		p		p		l		e				p		i		e		
0		1		2		3		4		5		6		7		8			

If you try to access an index beyond the end of a string, Python raises an `IndexError`:

```
>>> flavor[9]  
Traceback (most recent call last):  
  File "<pyshell#4>", line 1, in <module>  
    flavor[9]  
IndexError: string index out of range
```

The largest index in a string is always one less than the string's length. Since "apple pie" has a length of nine, the largest index allowed is 8.

Strings also support negative indices:

```
>>> flavor[-1]  
'e'
```

The last character in a string has index -1, which for "apple pie" is the letter e. The second-to-last character i has index -2, and so on.

The following figure shows the negative index for each character in the string "apple pie":

	a		p		p		l		e			p		i		e	
-9		-8		-7		-6		-5		-4		-3		-2		-1	

Just like positive indices, Python raises an `IndexError` if you try to access a negative index less than the index of the first character in the string:

```
>>> flavor[-10]
Traceback (most recent call last):
  File "<pysHELL#5>", line 1, in <module>
    flavor[-10]
IndexError: string index out of range
```

Negative indices may not seem useful at first, but sometimes they are a better choice than a positive index.

For example, suppose a string input by a user is assigned to the variable `user_input`. If you need to get the last character of the string, how do you know what index to use?

One way to get the last character of a string is to calculate the final index using `len()`:

```
final_index = len(user_input) - 1
last_character = user_input[final_index]
```

Getting the final character with the index `-1` takes less typing and doesn't require an intermediate step to calculate the final index:

```
last_character = user_input[-1]
```

String Slicing

Suppose you need the string containing just the first three letters of the string "apple pie". You could access each character by index and concatenate them, like this:

```
>>> first_three_letters = flavor[0] + flavor[1] + flavor[2]
>>> first_three_letters
'app'
```

If you need more than just the first few letters of a string, getting each character individually and concatenating them together is clumsy and long-winded. Fortunately, Python provides a way to do this with much less typing.

You can extract a portion of a string, called a **substring**, by inserting a colon between two index numbers inside of square brackets, like this:

```
>>> flavor = "apple pie"
>>> flavor[0:3]
'app'
```

`flavor[0:3]` returns the first three characters of the string assigned to `flavor`, starting with the character with index 0 and going up to, but not including, the character with index 3. The `[0:3]` part of `flavor[0:3]` is called a **slice**. In this case, it returns a slice of "apple pie". Yum!

String slices can be confusing because the substring returned by the slice includes the character whose index is the first number, but doesn't include the character whose index is the second number.

To remember how slicing works, you can think of a string as a sequence of square slots. The left and right boundary of each slot is numbered from zero up to the length of the string, and each slot is filled with a character in the string.

Here's what this looks like for the string "apple pie":

```
| a | p | p | l | e |     | p | i | e |
```

0	1	2	3	4	5	6	7	8	9
---	---	---	---	---	---	---	---	---	---

The slice `[x:y]` returns the substring between the boundaries `x` and `y`. So, for "apple pie", the slice `[0:3]` returns the string "app", and the slice `[3:9]` returns the string "le pie".

If you omit the first index in a slice, Python assumes you want to start at index 0:

```
>>> flavor[:5]
'apple'
```

The slice `[:5]` is equivalent to the slice `[0:5]`, so `flavor[:5]` returns the first five characters in the string "apple pie".

Similarly, if you omit the second index in the slice, Python assumes you want to return the substring that begins with the character whose index is the first number in the slice and ends with the last character in the string:

```
>>> flavor[5:]
' pie'
```

For "apple pie", the slice `[5:]` is equivalent to the slice `[5:9]`. Since the character with index 5 is a space, `flavor[5:9]` returns the substring that starts with the space and ends with the last letter, which is " pie".

If you omit both the first and second numbers in a slice, you get a string that starts with the character with index 0 and ends with the last character. In other words, omitting both numbers in a slice returns the entire string:

```
>>> flavor[:]
'apple pie'
```

It's important to note that, unlike string indexing, Python won't raise an `IndexError` when you try to slice between boundaries before or after

the beginning and ending boundaries of a string:

```
>>> flavor[:14]
'apple pie'
>>> flavor[13:15]
''
```

In this example, the first line gets the slice from the beginning of the string up to but not including the fourteenth character. The string assigned to `flavor` has length nine, so you might expect Python to throw an error. Instead, any non-existent indices are ignored and the entire string "apple pie" is returned.

The second shows what happens when you try to get a slice where the entire range is out of bounds. `flavor[13:15]` attempts to get the thirteenth and fourteenth characters, which don't exist. Instead of raising an error, the **empty string** "" is returned.

You can use negative numbers in slices. The rules for slices with negative numbers are exactly the same as slices with positive numbers. It helps to visualize the string as slots with the boundaries labeled by negative numbers:

	a		p		p		l		e				p		i		e	
-9	-8	-7	-6	-5	-4	-3	-2	-1										

Just like before, the slice `[x:y]` returns the substring between the boundaries `x` and `y`. For instance, the slice `[-9:-6]` returns the first three letters of the string "apple pie":

```
>>> flavor[-9:-6]
'app'
```

Notice, however, that the right-most boundary does not have a negative index. The logical choice for that boundary would seem to be the number 0, but that doesn't work:

```
>>> flavor[-9:0]
'',
```

Instead of returning the entire string, `[-9:0]` returns the **empty string** `""`. This is because the second number in a slice must correspond to a boundary that comes after the boundary corresponding to the first number, but both `-9` and `0` correspond to the left-most boundary in the figure.

If you need to include the final character of a string in your slice, you can omit the second number:

```
>>> flavor[-9:]
'apple pie'
```

Strings Are Immutable

To wrap this section up, let's discuss an important property of string objects. Strings are **immutable**, which means that you can't change them once you've created them. For instance, see what happens when you try to assign a new letter to one particular character of a string:

```
>>> word = "goal"
>>> word[0] = "f"
Traceback (most recent call last):
  File "<pyshell#16>", line 1, in <module>
    word[0] = "f"
TypeError: 'str' object does not support item assignment
```

Python throws a `TypeError` and tells you that `str` objects don't support item assignment.

Note

The term `str` is Python's internal name for the string data type.

If you want to alter a string, you must create an entirely new string. To change the string `"goal"` to the string `"foal"`, you can use a string

slice to concatenate the letter "f" with everything but the first letter of the word "goal":

```
>>> word = "goal"  
>>> word = "f" + word[1:]  
>>> word  
'foal'
```

First assign the string "goal" to the variable `word`. Then concatenate the slice `word[1:]`, which is the string "oal", with the letter "f" to get the string "foal". If you're getting a different result here, make sure you're including the : colon character as part of the string slice.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Create a string and print its length using the `len()` function.
2. Create two strings, concatenate them, and print the resulting string.
3. Create two strings and use concatenation to add a space in-between them. Then print the result.
4. Print the string "zing" by using slice notation on the string "bazinga" to specify the correct range of characters.

[Leave feedback on this section »](#)

4.3 Manipulate Strings With Methods

Strings come bundled with special functions called **string methods** that can be used to work with and manipulate strings. There are numerous string methods available, but we'll focus on some of the most commonly used ones.

In this section, you will learn how to:

- Convert a string to upper or lower case
- Remove whitespace from string
- Determine if a string begins and ends with certain characters

Let's go!

Converting String Case

To convert a string to all lower case letters, you use the string's `.lower()` method. This is done by tacking `.lower()` on to the end of the string itself:

```
>>> "Jean-luc Picard".lower()  
'jean-luc picard'
```

The dot (.) tells Python that what follows is the name of a method—the `lower()` method in this case.

Note

We will refer to the names of string methods with a dot at the beginning of them. So, for example, the `.lower()` method is written with a dot, instead of `lower()`.

The reason we do this is to make it easy to spot functions that are string methods, as opposed to built-in functions like `print()` and `type()`.

String methods don't just work on string literals. You can also use the `.lower()` method on a string assigned to a variable:

```
>>> name = "Jean-luc Picard"  
>>> name.lower()  
'jean-luc picard'
```

The opposite of the `.lower()` method is the `.upper()` method, which converts every character in a string to upper case:

```
>>> loud_voice = "Can you hear me yet?"  
>>> loud_voice.upper()  
'CAN YOU HEAR ME YET?'
```

Compare the `.upper()` and `.lower()` string methods to the general-purpose `len()` function you saw in the last section. Aside from the different results of these functions, the important distinction here is how they are used.

The `len()` function is a stand-alone function. If you want to determine the length of the `loud_voice` string, you call the `len()` function directly, like this:

```
>>> len(loud_voice)  
20
```

On the other hand, `.upper()` and `.lower()` must be used in conjunction with a string. They do not exist independently.

Removing Whitespace From a String

Whitespace is any character that is printed as blank space. This includes things like spaces and **line feeds**, which are special characters that move output to a new line.

Sometimes you need to remove whitespace from the beginning or end of a string. This is especially useful when working with strings that come from user input, where extra whitespace characters may have been introduced by accident.

There are three string methods that you can use to remove whitespace from a string:

1. `.rstrip()`
2. `.lstrip()`
3. `.strip()`

`.rstrip()` removes whitespace from the right side of a string:

```
>>> name = "Jean-luc Picard      "
>>> name
'Jean-luc Picard      '
>>> name.rstrip()
'Jean-luc Picard'
```

In this example, the string "Jean-luc Picard " has five trailing spaces. Python doesn't remove any trailing spaces in a string automatically when the string is assigned to a variable. The `.rstrip()` method removes trailing spaces from the right-hand side of the string and returns a new string "Jean-luc Picard", which no longer has the spaces at the end.

The `.lstrip()` method works just like `.rstrip()`, except that it removes whitespace from the left-hand side of the string:

```
>>> name = "      Jean-luc Picard"
>>> name
'      Jean-luc Picard'
>>> name.lstrip()
'Jean-luc Picard'
```

To remove whitespace from both the left and the right sides of the string at the same time, use the `.strip()` method:

```
>>> name = "      Jean-luc Picard      "
>>> name
'      Jean-luc Picard      '
>>> name.strip()
'Jean-luc Picard'
```

Note

None of the `.rstrip()`, `.lstrip()`, and `.strip()` methods remove whitespace from the middle of the string. In each of the previous examples the space between "Jean-luc" and "Picard" is always preserved.

Determine if a String Starts or Ends With a Particular String

When you work with text, sometimes you need to determine if a given string starts with or ends with certain characters. You can use two string methods to solve this problem: `.startswith()` and `.endswith()`.

Let's look at an example. Consider the string "Enterprise". Here's how you use `.startswith()` to determine if the string starts with the letters e and n:

```
>>> starship = "Enterprise"
>>> starship.startswith("en")
False
```

You must tell `.startswith()` what characters to search for by providing a string containing those characters. So, to determine if "Enterprise" starts with the letters e and n, you call `.startswith("en")`. This returns `False`. Why do you think that is?

If you guessed that `.startswith("en")` returns `False` because "Enterprise" starts with a capital E, you're absolutely right! The `.startswith()` method is **case-sensitive**. To get `.startswith()` to return `True`, you need to provide it with the string "En":

```
>>> starship.startswith("En")
True
```

The `.endswith()` method is used to determine if a string ends with certain characters:

```
>>> starship.endswith("rise")
True
```

Just like `.startswith()`, the `.endswith()` method is case-sensitive:

```
>>> starship.endswith("risE")
False
```

Note

The `True` and `False` values are not strings. They are a special kind of data type called a **Boolean value**. You will learn more about Boolean values in Chapter 8.

String Methods and Immutability

Recall from the previous section that strings are immutable—they can't be changed once they have been created. Most string methods that alter a string, like `.upper()` and `.lower()`, actually return copies of the original string with the appropriate modifications.

If you aren't careful, this can introduce subtle bugs into your program. Try this out in IDLE's interactive window:

```
>>> name = "Picard"
>>> name.upper()
'PICARD'
>>> name
'Picard'
```

When you call `name.upper()`, nothing about `name` actually changes. If you need to keep the result, you need to assign it to a variable:

```
>>> name = "Picard"
>>> name = name.upper()
>>> name
'PICARD'
```

`name.upper()` returns a new string "PICARD", which is re-assigned to the `name` variable. This **overrides** the original string "Picard" assigned to "name".

Use IDLE to Discover Additional String Methods

Strings have lots of methods associated to them. The methods introduced in this section barely scratch the surface. IDLE can help you

find new string methods. To see how, first assign a string literal to a variable in the interactive window:

```
>>> starship = "Enterprise"
```

Next, type `starship` followed by a period, but do not hit `Enter`. You should see the following in the interactive window:

```
>>> starship.
```

Now wait for a couple of seconds. IDLE displays a list of every string method that you can scroll through with the arrow keys.

A related shortcut in IDLE is the ability to fill in text automatically without having to type in long names by hitting `Tab`. For instance, if you only type in `starship.u` and then hit the `Tab` key, IDLE automatically fills in `starship.upper` because there is only one method belonging to `starship` that begins with a `u`.

This even works with variable names. Try typing in just the first few letters of `starship` and, assuming you don't have any other names already defined that share those first letters, IDLE completes the name `starship` for you when you hit the `Tab` key.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Write a script that converts the following strings to lowercase: "Animals", "Badger", "Honey Bee", "Honeybadger". Print each lowercase string on a separate line.
2. Repeat Exercise 1, but convert each string to uppercase instead of lowercase.
3. Write a script that removes whitespace from the following strings:

```
string1 = "      Filet Mignon"  
string2 = "Brisket      "
```

```
string3 = " Cheeseburger "
```

Print out the strings with the whitespace removed.

4. Write a script that prints out the result of `.startswith("be")` on each of the following strings:

```
string1 = "Becomes"
```

```
string2 = "becomes"
```

```
string3 = "BEAR"
```

```
string4 = " bEautiful"
```

5. Using the same four strings from Exercise 4, write a script that uses string methods to alter each string so that `.startswith("be")` returns True for all of them.

[Leave feedback on this section »](#)

4.4 Interact With User Input

Now that you've seen how to work with string methods, let's make things interactive. In this section, you will learn how to get some input from a user with the `input()` function. You'll write a program that asks a user to input some text and then display that text back to them in uppercase.

Enter the following into IDLE's interactive window:

```
>>> input()
```

When you press `Enter`, it looks like nothing happens. The cursor moves to a new line, but a new `>>>` doesn't appear. Python is waiting for you to enter something!

Go ahead and type some text and press `Enter`:

```
>>> input()
Hello there!
'Hello there!'
>>>
```

The text you entered is repeated on a new line with single quotes. That's because `input()` returns any text entered by the user as a string.

To make `input()` a bit more user friendly, you can give it a **prompt** to display to the user. The prompt is just a string that you put in between the parentheses of `input()`. It can be anything you want: a word, a symbol, a phrase—anything that is a valid Python string.

The `input()` function displays the prompt and waits for the user to type something on their keyboard. When the user hits `Enter`, `input()` returns their input as a string that can be assigned to a variable and used to do something in your program.

To see how `input()` works, save and run the following script:

```
prompt = "Hey, what's up? "
user_input = input(prompt)
print("You said:", user_input)
```

When you run this script, you'll see `Hey, what's up?` displayed in the interactive window with a blinking cursor.

The single space at the end of the string `"Hey, what's up? "` makes sure that when the user starts to type, the text is separated from the prompt with a space. When you type a response and press `Enter`, your response is assigned to the `user_input` variable.

Here's a sample run of the program:

```
Hey, what's up? Mind your own business.
```

```
You said: Mind your own business.
```

Once you have input from a user, you can do something with it. For example, the following script takes user input and converts it to uppercase with `.upper()` and prints the result:

```
response = input("What should I shout? ")
shouted_response = response.upper()
print("Well, if you insist...", shouted_response)
```

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Write a script that takes input from the user and displays that input back.
2. Write a script that takes input from the user and displays the input in lowercase.
3. Write a script that takes input from the user and displays the number of characters inputted.

[Leave feedback on this section »](#)

4.5 Challenge: Pick Apart Your User's Input

Write a script named `first_letter.py` that first prompts the user for input by using the string "Tell me your password:" The script should then determine the first letter of the user's input, convert that letter to upper-case, and display it back.

For example, if the user input is "no" then the program should respond like this:

```
The first letter you entered was: N
```

For now, it's okay if your program crashes when the user enters nothing as input—that is, they just hit `Enter` instead of typing something in. You'll learn about a couple of ways you can deal with this situation in an upcoming chapter.

You can find the solutions to this code challenge and many other bonus resources online at realpython.com/python-basics/resources.

[Leave feedback on this section »](#)

4.6 Working With Strings and Numbers

When you get user input using the `input()` function, the result is always a string. There are many other times when input is given to a program as a string. Sometimes those strings contain numbers that need to be fed into calculations.

In this section you will learn how to deal with strings of numbers. You will see how arithmetic operations work on strings, and how they often lead to surprising results. You will also learn how to convert between strings and number types.

Strings and Arithmetic Operators

You've seen that string objects can hold many types of characters, including numbers. However, don't confuse numerals in a string with actual numbers. For instance, try this bit of code out in IDLE's interactive window:

```
>>> num = "2"  
>>> num + num  
'22'
```

The `+` operator concatenates two string together. So, the result of `"2" + "2"` is `"22"`, not `"4"`.

Strings can be “multiplied” by a number as long as that number is an integer, or whole number. Type the following into the interactive window:

```
>>> num = "12"  
>>> num * 3
```

```
'121212'
```

`num * 3` concatenates the string "12" with itself three times and returns the string "121212". To compare this operation to arithmetic with numbers, notice that `"12" * 3 = "12" + "12" + "12"`. In other words, multiplying a string by an integer `n` concatenates that string with itself `n` times.

The number on the right-hand side of the expression `num * 3` can be moved to the left, and the result is unchanged:

```
>>> 3 * num  
'121212'
```

What do you think happens if you use the `*` operator between two strings? Type `"12" * "3"` in the interactive window and press `Enter`:

```
>>> "12" * "3"  
Traceback (most recent call last):  
  File "<stdin>", line 1, in <module>  
TypeError: can't multiply sequence by non-int of type 'str'
```

Python raises a `TypeError` and tells you that you can't multiply a sequence by a non-integer. When the `*` operator is used with a string on either the left or the right side, it always expects an integer on the other side.

Note

A **sequence** is any Python object that supports accessing elements by index. Strings are sequences. You will learn about other sequence types in Chapter 9.

What do you think happens when you try to add a string and a number?

```
>>> "3" + 3  
Traceback (most recent call last):
```

```
File "<stdin>", line 1, in <module>
TypeError: can only concatenate str (not "int") to str
```

Again, Python throws a `TypeError` because the `+` operator expects both things on either side of it to be of the same type. If any one of the objects on either side of `+` is a string, Python tries to perform string concatenation. Addition will only be performed if both objects are numbers. So, to add `"3" + 3` and get `6`, you must first convert the string `"3"` to a number.

Converting Strings to Numbers

The `TypeError` errors you saw in the previous section highlight a common problem encountered when working with user input: type mismatches when trying to use the input in an operation that requires a number and not a string.

Let's look at an example. Save and run the following script.

```
num = input("Enter a number to be doubled: ")
doubled_num = num * 2
print(doubled_num)
```

When you enter a number, such as `2`, you expect the output to be `4`, but in this case, you get `22`. Remember, `input()` always returns a string, so if you input `2`, then `num` is assigned the string `"2"`, not the integer `2`. Therefore, the expression `num * 2` returns the string `"2"` concatenated with itself, which is `"22"`.

To perform arithmetic on numbers that are contained in a string, you must first convert them from a string type to a number type. There are two ways to do this: `int()` and `float()`.

`int()` stands for **integer** and converts objects into whole numbers, while `float()` stands for **floating-point number** and converts objects into numbers with decimal points. Here's what using them looks like in the interactive window:

```
>>> int("12")
12

>>> float("12")
12.0
```

Notice how `float()` adds a decimal point to the number. Floating-point numbers always have at least one decimal place of precision. For this reason, you can't change a string that looks like a floating-point number into an integer because you would lose everything after the decimal point:

```
>>> int("12.0")
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
ValueError: invalid literal for int() with base 10: '12.0'
```

Even though the extra `0` after the decimal place doesn't add any value to the number, Python won't change `12.0` into `12` because it would result in the loss of precision.

Let's revisit the script from the beginning of this section and see how to fix it. Here's the script again:

```
num = input("Enter a number to be doubled: ")
doubled_num = num * 2
print(doubled_num)
```

The issue lies in the line `doubled_num = num * 2` because `num` references a string and `2` is an integer. You can fix the problem by wrapping `num` with either `int()` or `float()`. Since the prompts asks the user to input a number, and not specifically an integer, let's convert `num` to a floating-point number:

```
num = input("Enter a number to be doubled: ")
doubled_num = float(num) * 2
print(doubled_num)
```

Now when you run this script and input 2, you get 4.0 as expected. Try it out!

Converting Numbers to Strings

Sometimes you need to convert a number to a string. You might do this, for example, if you need to build a string from some pre-existing variables that are assigned to numeric values.

As you've already seen, the following produces a `TypeError`:

```
>>> num_pancakes = 10
>>> "I am going to eat " + num_pancakes + " pancakes."
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
TypeError: can only concatenate str (not "int") to str
```

Since `num_pancakes` is a number, Python can't concatenate it with the string "I'm going to eat". To build the string, you need to convert `num_pancakes` to a string using `str()`:

```
>>> num_pancakes = 10
>>> "I am going to eat " + str(num_pancakes) + " pancakes."
'I am going to eat 10 pancakes.'
```

You can also call `str()` on a number literal:

```
>>> "I am going to eat " + str(10) + " pancakes."
'I am going to eat 10 pancakes.'
```

`str()` can even handle arithmetic expressions:

```
>>> total_pancakes = 10
>>> pancakes_eaten = 5
>>> "Only " + str(total_pancakes - pancakes_eaten) + " pancakes left."
'Only 5 pancakes left.'
```

You're not limited to numbers when using `str()`. You can pass it all sorts of objects to get their string representations:

```
>>> str(print)
'<built-in function print>'

>>> str(int)
"<class 'int'>"

>>> str(float)
"<class 'float'>"
```

These examples may not seem very useful, but they illustrate how flexible `str()` is.

In the next section, you'll learn how to format strings neatly to display values in a nice, readable manner. Before you move on, though, check your understanding with the following review exercises.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Create a string containing an integer, then convert that string into an actual integer object using `int()`. Test that your new object is a number by multiplying it by another number and displaying the result.
2. Repeat the previous exercise, but use a floating-point number and `float()`.
3. Create a string object and an integer object, then display them side-by-side with a single `print` statement by using the `str()` function.
4. Write a script that gets two numbers from the user using the `input()` function twice, multiplies the numbers together, and displays the result. If the user enters 2 and 4, your program should print the following text:
`The product of 2 and 4 is 8.0.`

[Leave feedback on this section »](#)

4.7 Streamline Your Print Statements

Suppose you have a string `name = "Zaphod"` and two integers `heads = 2` and `arms = 3`. You want to display them in the following line: `Zaphod has 2 heads and 3 arms.` This is called **string interpolation**, which is just a fancy way of saying that you want to insert some variables into specific locations in a string.

You've already seen two ways of doing this. The first involves using commas to insert spaces between each part of the string inside of a `print()` function:

```
print(name, "has", str(heads), "heads and", str(arms), "arms")
```

Another way to do this is by concatenating the strings with the `+` operator:

```
print(name + " has " + str(heads) + " heads and " + str(arms) + " arms")
```

Both techniques produce code that can be hard to read. Trying to keep track of what goes inside or outside of the quotes can be tough. Fortunately, there's a third way of combining strings: **formatted string literals**, more commonly known as f-strings.

The easiest way to understand f-strings is to see them in action. Here's what the above string looks like when written as an f-string:

```
>>> f"{name} has {heads} heads and {arms} arms"  
'Zaphod has 2 heads and 3 arms'
```

There are two important things to notice about the above examples:

1. The string literal starts with the letter `f` before the opening quotation mark
2. Variable names surrounded by curly braces (`{` and `}`) are replaced with their corresponding values without using `str()`

You can also insert Python expressions in between the curly braces. The expressions are replaced with their result in the string:

```
>>> n = 3
>>> m = 4
>>> f"{n} times {m} is {n*m}"
'3 times 4 is 12'
```

It is a good idea to keep any expressions used in an f-string as simple as possible. Packing in a bunch of complicated expressions into a string literal can result in code that is difficult to read and difficult to maintain.

f-strings are only available in Python version 3.6 and above. In earlier versions of Python, the `.format()` method can be used to get the same results. Returning to the Zaphod example, you can use `.format()` method to format the string like this:

```
>>> "{} has {} heads and {} arms".format(name, heads, arms)
'Zaphod has 2 heads and 3 arms'
```

f-strings are shorter, and sometimes more readable, than using `.format()`. You will see f-strings used throughout this book.

For an in-depth guide to f-strings and comparisons to other string formatting techniques, check out the [Python 3's f-Strings: An Improved String Formatting Syntax \(Guide\)](#) on realpython.com

Note

There is also another way to print formatted strings: using the `%` operator. You might see this in code that you find elsewhere, and you can [read about how it works here](#) if you're curious.

Keep in mind that this style has been phased out entirely in Python 3. Just be aware that it exists and you may see it in legacy Python code bases.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Create a `float` object named `weight` with the value `0.2`, and create a string object named `animal` with the value `"newt"`. Then use these objects to print the following string using only string concatenation:
0.2 kg is the weight of the newt.
2. Display the same string by using the `.format()` method and empty `{}` place-holders.
3. Display the same string using an f-string.

[Leave feedback on this section »](#)

4.8 Find a String in a String

One of the most useful string methods is `.find()`. As its name implies, you can use this method to find the location of one string in another string—commonly referred to as a **substring**.

To use `.find()`, tack it to the end of a variable or a string literal and pass the string you want to find in between the parentheses:

```
>>> phrase = "the surprise is in here somewhere"
>>> phrase.find("surprise")
4
```

The value that `.find()` returns is the index of the first occurrence of the string you pass to it. In this case, "surprise" starts at the fifth character of the string "the surprise is in here somewhere" which has index 4 because counting starts at 0.

If `.find()` doesn't find the desired substring, it will return `-1` instead:

```
>>> phrase = "the surprise is in here somewhere"  
>>> phrase.find("eyjafjallajökull")  
-1
```

You can call string methods on a string literal directly, so in this case, you don't need to create a new string:

```
>>> "the surprise is in here somewhere".find("surprise")  
4
```

Keep in mind that this matching is done exactly, character by character, and is case-sensitive. For example, if you try to find "SURPRISE", the `.find()` method returns -1:

```
>>> "the surprise is in here somewhere".find("SURPRISE")  
-1
```

If a substring appears more than once in a string, `.find()` only returns the index of the first appearance, starting from the beginning of the string:

```
>>> "I put a string in your string".find("string")  
8
```

There are two instances of the "string" in "I put a string in your string". The first starts at index 8, and the second at index 23. `.find()` returns 8, which is the index of the first instance of "string".

The `.find()` method only accepts a string as its input. If you want to find an integer in a string, you need to pass the integer to `.find()` as a string. If you do pass something other than a string to `.find()`, Python raises a `TypeError`:

```
>>> "My number is 555-555-5555".find(5)  
Traceback (most recent call last):  
  File "<stdin>", line 1, in <module>  
TypeError: must be str, not int
```

```
>>> "My number is 555-555-5555".find("5")
```

13

Sometimes you need to find all occurrences of a particular substring and replace them with a different string. Since `.find()` only returns the index of the first occurrence of a substring, you can't easily use it to perform this operation. Fortunately, string objects have a `.replace()` method that replaces each instance of a substring with another string.

Just like `.find()`, you tack `.replace()` on to the end of a variable or string literal. In this case, though, you need to put two strings inside of the parentheses in `.replace()` and separate them with a comma. The first string is the substring to find, and the second string is the string to replace each occurrence of the substring with.

For example, the following code shows how to replace each occurrence of "the truth" in the string "I'm telling you the truth; nothing but the truth" with the string "lies":

```
>>> my_story = "I'm telling you the truth; nothing but the truth!"  
>>> my_story.replace("the truth", "lies")  
"I'm telling you lies; nothing but lies!"
```

Since strings are immutable objects, `.replace()` doesn't alter `my_story`. If you immediately type `my_story` into the interactive window after running the above example, you'll see the original string, unaltered:

```
>>> my_story  
"I'm telling you the truth; nothing but the truth!"
```

To change the value of `my_story`, you need to reassign to it the new value returned by `.replace()`:

```
>>> my_story = my_story.replace("the truth", "lies")  
>>> my_story  
"I'm telling you lies; nothing but lies!"
```

`.replace()` can only replace one substring at a time, so if you want to replace multiple substrings in a string you need to use `.replace()` mul-

tiple times:

```
>>> text = "some of the stuff"  
>>> new_text = text.replace("some of", "all")  
>>> new_text = new_text.replace("stuff", "things")  
>>> new_text  
'all the things'
```

You'll have some fun with `.replace()` in the challenge in the next section.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. In one line of code, display the result of trying to `.find()` the substring "a" in the string "AAA". The result should be -1.
2. Replace every occurrence of the character "s" with "x" in the string "Somebody said something to Samantha.".
3. Write and test a script that accepts user input using the `input()` function and displays the result of trying to `.find()` a particular letter in that input.

[Leave feedback on this section »](#)

4.9 Challenge: Turn Your User Into a L33t H4xor

Write a script called `translate.py` that asks the user for some input with the following prompt: `Enter some text::`. Then use the `.replace()` method to convert the text entered by the user into “leetspeak” by making the following changes to lower-case letters:

- The letter a becomes 4
- The letter b becomes 8

- The letter e becomes 3
- The letter l becomes 1
- The letter o becomes 0
- The letter s becomes 5
- The letter t becomes 7

Your program should then display the resulting string as output. Below is a sample run of the program:

```
Enter some text: I like to eat eggs and spam.  
I lik3 70 347 3gg5 4nd 5p4m.
```

You can find the solutions to this code challenge and many other bonus resources online at realpython.com/python-basics/resources.

[Leave feedback on this section »](#)

4.10 Summary and Additional Resources

In this chapter, you learned the ins and outs of Python string objects. You learned how to access different characters in a string using subscripts and slices, as well as how to determine the length of a string with `len()`.

Strings come with numerous methods. The `.upper()` and `.lower()` methods convert all characters of a string to upper or lower case, respectively. The `.rstrip()`, `.lstrip()`, and `strip()` methods remove whitespace from strings, and the `.startswith()` and `.endswith()` methods will tell you if a string starts or ends with a given substring.

You also saw how to capture input from a user as a string using the `input()` function, and how to convert that input to a number using `int()` and `float()`. To convert numbers, and other objects, to strings, you use `str()`.

Finally, you saw how the `.find()` and `.replace()` methods are used to find the location of a substring and replace a substring with a new string.

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-4

Additional Resources

To learn more, check out the following resources:

- [Python String Formatting Best Practices](#)
- [Splitting, Concatenating, and Joining Strings in Python](#)
- [Recommended resources on realpython.com](#)

[Leave feedback on this section »](#)

Chapter 5

Numbers and Math

You don't need to be a math whiz to program well. The truth is, few programmers need to know more than basic algebra.

Of course, how much math you need to know depends on the application you are working on. In general, the level of math required to successfully work as a programmer is less than you might expect.

Although math and computer programming aren't as correlated as some people might believe, numbers are an integral part of any programming language and Python is no exception.

In this chapter, you will learn how to:

- Work with Python's three built-in number types: integer, floating-point, and complex numbers
- Round numbers to a given number of decimal places
- Format and display numbers in strings

Let's get started!

[Leave feedback on this section »](#)

5.1 Integers and Floating-Point Numbers

Python has three built-in number data types: integers, floating-point numbers, and complex numbers. In this section, you'll learn about integers and floating-point numbers, which are the two most commonly used number types. You'll learn about complex numbers in section 5.6.

Integers

An **integer** is a whole number with no decimal places. For example, `1` is an integer, but `1.0` isn't. The name for the integer data type is `int`, which you can see with the `type()` function:

```
>>> type(1)
<class 'int'>
```

You can create an integer by simply typing the number explicitly or using the `int()` function. In Chapter 4, you learned how to convert a string containing an integer to a number using `int()`. For example, the following converts the string "25" to the integer 25:

```
>>> int("25")
25
```

An **integer literal** is an integer value that is written explicitly in your code, just like a string literal is a string that is written explicitly in your code. For example, `1` is an integer literal, but `int("1")` isn't.

Integer literals can be written in two different ways:

```
>>> 1000000
1000000

>>> 1_000_000
1000000
```

The first example is straightforward. Just type a 1 followed by six zeros. The downside to this notation is that large numbers can be difficult to read.

When you write large numbers by hand, you probably group digits into groups of three, separated by a comma. 1,000,000 is a lot easier to read than 1000000.

In Python, you can't use commas to group digits in integer literals, but you can use an underscore (_). The value 1_000_000 expresses one million in a more readable manner.

There is no limit to how large an integer can be, which might be surprising considering computers have finite memory. Try typing the largest number you can think of into IDLE's interactive window. Python can handle it with no problem!

Floating-Point Numbers

A **floating-point number**, or **float** for short, is a number with a decimal place. 1.0 is a floating-point number, as is -2.75. The name of a floating-point data type is **float**:

```
>>> type(1.0)
<class 'float'>
```

Floats can be created by typing a number directly into your code, or by using the `float()` function. Like `int()`, `float()` can be used to convert a string containing a number to a floating-point number:

```
>>> float("1.25")
1.25
```

A **floating-point literal** is a floating-point value that is written explicitly in your code. 1.25 is a floating-point literal, while `float("1.25")` is not.

Floating-point literals can be created in three different ways. Each of the following creates a floating-point literal with a value of one mil-

lion:

```
>>> 1000000.0  
1000000.0  
  
>>> 1_000_000.0  
1000000.0  
  
>>> 1e6  
1000000.0
```

The first two ways are similar to the two methods for creating integer literals that you saw earlier. The second method, which uses underscores to separate digits into groups of three, is useful for creating float literals with lots of digits.

For really large numbers, you can use **E-notation**. The third method in the previous example uses E-notation to create a float literal.

To write a float literal in E-notation, type a number followed by the letter e and then another number. Python takes the number to the left of the e and multiplies by 10 raised to the power of the number after the e. So 1e6 is equivalent to 1×10^6 .

Note

E-notation is short for **exponential notation**, and is the more common name for how many calculators and programming languages display large numbers.

Python also uses E-notation to display large floating point numbers:

```
>>> 2000000000000000000.0  
2e+17
```

The float 2000000000000000000.0 gets displayed as 2e+17. The + sign indicates that the exponent 17 is a positive number. You can also use negative numbers as the exponent:

```
>>> 1e-4  
0.0001
```

The literal `1e-4` is interpreted as `10` raised to the power `-4`, which is `1/10000` or, equivalently, `0.0001`.

Unlike integers, floats do have a maximum size. The maximum floating-point number depends on your system, but something like `2e400` ought to be well beyond most machines' capabilities. `2e400` is 2×10^{400} , which is far more than the [total number of atoms in the universe!](#)

When you reach the maximum floating-point number, Python returns a special float value `inf`:

```
>>> 2e400  
inf
```

`inf` stands for infinity, and it just means that the number you've tried to create is beyond the maximum floating-point value allowed on your computer. The type of `inf` is still `float`:

```
>>> n = 2e400  
>>> n  
inf  
>>> type(n)  
<class 'float'>
```

There is also `-inf` which stands for negative infinity, and represents a negative floating-point number that is beyond the minimum floating-point number allowed on your computer:

```
>>> -2e400  
-inf
```

You probably won't come across `inf` and `-inf` often as a programmer, unless you regularly work with extremely large numbers.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Write a script that creates the two variables, `num1` and `num2`. Both `num1` and `num2` should be assigned the integer literal 25,000,000, one written with underscored and one without. Print `num1` and `num2` on two separate lines.
2. Write a script that assigns the floating-point literal 175000.0 to the variable `num` using exponential notation, and then prints `num` in the interactive window.
3. In IDLE’s interactive window, try and find the smallest exponent `N` so that `2e<N>`, where `<N>` is replaced with your number, returns `inf`.

[Leave feedback on this section »](#)

5.2 Arithmetic Operators and Expressions

In this section, you’ll learn how to do basic arithmetic with numbers in Python, such as addition, subtraction, multiplication, and division. Along the way, you’ll learn some conventions for writing mathematical expressions in code.

Addition

Addition is performed with the `+` operator:

```
>>> 1 + 2  
3
```

The two numbers on either side of the `+` operator are called **operands**. In the previous example, both operands are integers, but operands do not need to be the same type. You can add an `int` to a `float` with no problem:

```
>>> 1.0 + 2  
3.0
```

Notice that the result of `1.0 + 2` is `3.0`, which is a `float`. Any time a `float` is added to a number, the result is another `float`. Adding two integers together always results in an `int`.

Note

PEP 8 recommends separating both operands from an operator with a space.

Python can evaluate `1+1` just fine, but `1 + 1` is the preferred format because it's generally considered easier to read. This rule of thumb applies to all of the operators in this section.

Subtraction

To subtract two numbers, just put a `-` in between them:

```
>>> 1 - 1  
0  
  
>>> 5.0 - 3  
2.0
```

Just like adding two integers, subtracting two integers always results in an `int`. Whenever one of the operands is a `float`, the result is also a `float`.

The `-` operator is also used to denote negative numbers:

```
>>> -3  
-3
```

You can subtract a negative number from another number, but as you can see below, this can sometimes look confusing:

```
>>> 1 - -3
```

```
4
```

```
>>> 1 --3
```

```
4
```

```
>>> 1- -3
```

```
4
```

```
>>> 1--3
```

```
4
```

Of the four examples above, the first is the most PEP 8 compliant. That said, you can surround `-3` with parentheses to make it even clearer that the second `-` is modifying `3`:

```
>>> 1 - (-3)
```

```
4
```

Using parentheses is a good idea because it makes the code more explicit. Computers execute code, but humans read code. Anything you can do to make your code easier to read and understand is a good thing.

Multiplication

To multiply two numbers, use the `*` operator:

```
>>> 3 * 3
```

```
9
```

```
>>> 2 * 8.0
```

```
16.0
```

The type of number you get from multiplication follows the same rules as addition and subtraction. Multiplying two integers results in an `int`, and multiplying a number with a `float` results in a `float`.

Division

The `/` operator is used to divide two numbers:

```
>>> 9 / 3
```

```
3.0
```

```
>>> 5.0 / 2
```

```
2.5
```

Unlike addition, subtraction, and multiplication, division with the `/` operator always returns a `float`. If you want to make sure that you get an integer after dividing two numbers, you can use `int()` to convert the result:

```
>>> int(9 / 3)
```

```
3
```

Keep in mind that `int()` discards any fractional part of the number:

```
>>> int(5.0 / 2)
```

```
2
```

`5.0 / 2` returns the float `2.5`, and `int(2.5)` returns the integer `2` with the `.5` part removed.

Integer Division

If writing `int(5.0 / 2)` seems a little long-winded to you, Python provides a second division operator, `//`, called the **integer division** operator:

```
>>> 9 // 3
```

```
3
```

```
>>> 5.0 // 2
```

```
2.0
```

```
>>> -3 // 2  
-2
```

The `//` operator first divides the number on the left by the number on the right and then rounds down to an integer. This might not give the value you expect when one of the numbers is negative.

For example, `-3 // 2` returns `-2`. First, `-3` is divided by `2` to get `-1.5`. Then `-1.5` is rounded down to `-2`. On the other hand, `3 // 2` returns `1`.

Another thing the above example illustrates is that `//` returns a floating-point number if one of the operands is a float. This is why `9 // 3` returns the integer `3` and `5.0 // 2` returns the float `2.0`.

Let's see what happens when you try to divide a number by `0`:

```
>>> 1 / 0  
Traceback (most recent call last):  
  File "<stdin>", line 1, in <module>  
ZeroDivisionError: division by zero
```

Python gives you a `ZeroDivisionError`, letting you know that you just tried to break a fundamental rule of the universe.

Exponents

You can raise a number to a power using the `**` operator:

```
>>> 2 ** 2  
4  
  
>>> 2 ** 3  
8  
  
>>> 2 ** 4  
16
```

Exponents don't have to be integers. They can also be floats:

```
>>> 3 ** 1.5  
5.196152422706632
```

```
>>> 9 ** 0.5  
3.0
```

Raising a number to the power of 0.5 is the same as taking the square root, but notice that even though the square root of 9 is an integer, Python returns the float 3.0.

For positive operands, the `**` operator returns an integer if both operands are integers, and a float if any one of the operands is a floating-point number.

You can also raise numbers to negative powers:

```
>>> 2 ** -1  
0.5
```

```
>>> 2 ** -2  
0.25
```

Raising a number to a negative power is the same as dividing 1 by the number raised to the positive power. So, `2 ** -1` is the same as `1 / (2 ** 1)`, which is the same as `1 / 2`, or 0.5. Similarly `2 ** -2` is the same as `1 / (2 ** 2)`, which is the same as `1 / 4`, or 0.25.

The Modulus Operator

The `%` operator, or the **modulus**, returns the remainder of dividing the left operand by the right operand:

```
>>> 5 % 3
```

```
2
```

```
>>> 20 % 7
```

```
6
```

```
>>> 16 % 8  
0
```

3 divides 5 once with a remainder of 2, so $5 \% 3$ is 2. Similarly, 7 divides 20 twice with a remainder of 6.

In the last example, 16 is divisible by 8, so $16 \% 8$ is 0. Any time the number to the left of `%` is divisible by the number to the right, the result is 0.

One of the most common uses of `%` is to determine whether or not one number is divisible by another. For example, a number n is even if and only if $n \% 2$ is 0.

What do you think $1 \% 0$ returns? Let's try it out:

```
>>> 1 % 0  
Traceback (most recent call last):  
  File "<stdin>", line 1, in <module>  
ZeroDivisionError: integer division or modulo by zero
```

This makes sense because $1 \% 0$ is the remainder of dividing 1 by 0. But you can't divide 1 by 0, so Python raises a `ZeroDivisionError`.

Note

When you work in IDLE's interactive window, errors like `ZeroDivisionError` don't cause much of a problem. The error is displayed and a new prompt pops up allowing you to continue writing code.

However, whenever Python encounters an error while running a script, execution stops. The program is said to have **crashed**. In Chapter 8, you'll learn how to handle errors so that your programs don't crash unexpectedly.

Things get a little tricky when you use the `%` operator with negative numbers:

```
>>> 5 % -3
```

```
-1
```

```
>>> -5 % 3
```

```
1
```

```
>>> -5 % -3
```

```
-2
```

These potentially shocking results are really quite well defined. To calculate the remainder r of dividing a number x by a number y , Python uses the equation $r = x - (y * (x // y))$.

For example, to find $5 \% -3$, first find $(5 // -3)$. Since $5 / -3$ is about -1.67 , $5 // -3$ is -2 . Now multiply that by -3 to get 6 . Finally, subtract 6 from 5 to get -1 .

Arithmetic Expressions

You can combine operators to form complex expressions. An **expression** is a combination of numbers, operators, and parentheses that Python can compute, or **evaluate**, to return a value.

Here are some examples of arithmetic expressions:

```
>>> 2**3 - 1
```

```
5
```

```
>>> 4/2 + 2**3
```

```
10.0
```

```
>>> -1 + (-3**2 + 4)
```

```
-3
```

The rules for evaluating expressions work are the same as in everyday arithmetic. In school, you probably learned these rules under the name “order of operations.”

The `*`, `/`, `//`, and `%` operators all have equal **precedence**, or priority, in an expression, and each of these has a higher precedence than the `+` and `-` operators. This is why `2*3 - 1` returns `5` and not `4`. `2*3` is evaluated first, because `*` has higher precedence than the `-` operator.

You may notice that the expressions in the previous example do not follow the rule for putting a space on either side of all of the operators. PEP 8 says the following about whitespace in complex expressions:

“If operators with different priorities are used, consider adding whitespace around the operators with the lowest priority(ies). Use your own judgment; however, never use more than one space, and always have the same amount of whitespace on both sides of a binary operator.”

— PEP 8, Other Recommendations

[Leave feedback on this section »](#)

5.3 Challenge: Perform Calculations on User Input

Write a script called `exponent.py` that receives two numbers from the user and displays the first number raised to the power of the second number.

A sample run of the program should look like this (with example input that has been provided by the user included below):

```
Enter a base: 1.2
Enter an exponent: 3
1.2 to the power of 3 = 1.7279999999999998
```

Keep the following in mind:

1. Before you can do anything with the user's input, you will have to assign both calls to `input()` to new variables.
2. The `input()` function returns a string, so you'll need to convert the user's input into numbers in order to do arithmetic.
3. You can use an f-string to print the result.
4. You can assume that the user will enter actual numbers as input.

You can find the solutions to this code challenge and many other bonus resources online at realpython.com/python-basics/resources.

[Leave feedback on this section »](#)

5.4 Make Python Lie to You

What do you think $0.1 + 0.2$ is? The answer is 0.3 , right? Let's see what Python has to say about it. Try this out in the interactive window:

```
>>> 0.1 + 0.2  
0.3000000000000004
```

Well, that's... *almost* right! What in the heck is going on here? Is this a bug in Python?

No, it isn't a bug! It's a **floating-point representation error**, and it has nothing to do with Python. It's related to the way floating-point numbers are stored in a computer's memory.

The number 0.1 can be represented as the fraction $1/10$. Both the number 0.1 and its fraction $1/10$ are **decimal representations**, or **base 10 representations**. Computers, however, store floating-point numbers in base 2 representation, more commonly called **binary representation**.

When represented in binary, something familiar yet possibly unexpected happens to the decimal number 0.1 . The fraction $1/3$ has no finite decimal representation. That is, $1/3 = 0.3333\dots$ with infinitely many 3 's after the decimal point. The same thing happens to the frac-

tion $1/10$ in binary.

The binary representation of $1/10$ is the following infinitely repeating fraction:

```
0.00011001100110011001100110011...
```

Computers have finite memory, so the number 0.1 must be stored as an approximation and not as its true value. The approximation that gets stored is slightly higher than the actual value, and looks like this:

```
0.100000000000000055511151231257827021181583404541015625
```

You may have noticed, however, that when asked to print 0.1 , Python prints 0.1 and not the approximated value above:

```
>>> 0.1  
0.1
```

Python doesn't just chop off the digits in the binary representation for 0.1 . What actually happens is a little more subtle.

Because the approximation of 0.1 in binary is just that—an approximation—it is entirely possible that more than one decimal number have the same binary approximation.

For example, the numbers 0.1 and 0.1000000000000001 both have the same binary approximation. Python prints out the shortest decimal number that shares the approximation.

This explains why, in the first example of this section, $0.1 + 0.2$ does not equal 0.3 . Python adds together the binary approximations for 0.1 and 0.2 , which gives a number which is *not* the binary approximation for 0.3 .

If all this is starting to make your head spin, don't worry! Unless you are writing programs for finance or scientific computing, you don't need to worry about the imprecision of floating-point arithmetic.

[Leave feedback on this section »](#)

5.5 Math Functions and Number Methods

Python has a few built-in functions you can use to work with numbers. In this section, you'll learn about three of the most common ones:

1. `round()`, for rounding numbers to some number of decimal places
2. `abs()`, for getting the absolute value of a number
3. `pow()`, for raising a number to some power

You'll also learn about a method that floating-point numbers have to check whether or not they have an integer value.

Let's go!

The `round()` function

You can use `round()` to round a number to the nearest integer:

```
>>> round(2.3)
```

```
2
```

```
>>> round(2.7)
```

```
3
```

`round()` has some unexpected behavior when the number ends in `.5`:

```
>>> round(2.5)
```

```
2
```

```
>>> round(3.5)
```

```
4
```

`2.5` gets rounded down to `2` and `3.5` is rounded up to `4`. Most people expect a number that ends in `.5` to get rounded up, so let's take a closer look at what's going on here.

Python 3 rounds numbers according to a strategy called **rounding**

ties to even. A **tie** is any number whose last digit is a five. 2.5 and 3.1415 are ties, but 1.37 is not.

When you round ties to even, you first look at the digit one decimal place to the left of the last digit in the tie. If that digit is even, you round down. If the digit is odd, you round up. That's why 2.5 rounds down to 2 and 3.5 round up to 4.

Note

Rounding ties to even is the rounding strategy recommended for floating-point numbers by the [IEEE](#) (Institute of Electrical and Electronics Engineers) because it helps limit the impact rounding has on operations involving lots of numbers.

The IEEE maintains a standard called [IEEE 754](#) for how floating-point numbers are dealt with on a computer. It was published in 1985 and is still commonly used by hardware manufacturers today.

You can round a number to a given number of decimal places by passing a second argument to `round()`:

```
>>> round(3.14159, 3)
```

```
3.142
```

```
>>> round(2.71828, 2)
```

```
2.72
```

The number 3.14159 is rounded to 3 decimal places to get 3.142, and the number 2.71828 is rounded to 2 decimal places to get 2.72.

The second argument of `round()` must be an integer. If it isn't, Python raises a `TypeError`:

```
>>> round(2.65, 1.4)
```

```
Traceback (most recent call last):
```

```
  File "<pyshell#0>", line 1, in <module>
```

```
round(2.65, 1.4)
TypeError: 'float' object cannot be interpreted as an integer
```

Sometimes `round()` doesn't get the answer quite right:

```
>>> # Expected value: 2.68
>>> round(2.675, 2)
2.67
```

2.675 is a tie because it lies exactly halfway between the numbers 2.67 and 2.68. Since Python rounds ties to the nearest even number, you would expect `round(2.675, 2)` to return 2.68, but it returns 2.67 instead. This error is a result of floating-point representation error, and isn't a bug in the `round()` function.

Dealing with floating-point numbers can be frustrating, but this frustration isn't specific to Python. All languages that implement the IEEE floating-point standard have the same issues, including C/C++, Java, and JavaScript.

In most cases, though, the little errors encountered with floating-point numbers are negligible, and the results of `round()` are perfectly useful.

The `abs()` Function

The **absolute value** of a number n is just n if n is positive, and $-n$ if n is negative. For example, the absolute value of 3 is 3, while the absolute value of -5 is 5.

To get the absolute value of a number in Python, you use the `abs()` function:

```
>>> abs(3)
3

>>> abs(-5.0)
5.0
```

`abs()` always returns a positive number of the same type as its argument. That is, the absolute value of an integer is always a positive integer, and the absolute value of a float is always a positive float.

The `pow()` Function

In section 5.2, you learned how to raise a number to a power using the `**` operator. You can also use the `pow()` function. `pow()` takes two arguments. The first is the **base**, that is the number to be raised to a power, and the second argument is the **exponent**.

For example, the following uses `pow()` to raise 2 to the exponent 3:

```
>>> pow(2, 3)
8
```

Just like `**`, the exponent in `pow()` can be negative:

```
>>> pow(2, -2)
0.25
```

So, what's the difference between `**` and `pow()`? The `pow()` function accepts an optional third argument that computes the first number raised to the power of the second number and then takes the modulo with respect to the third number.

In other words, `pow(x, y, z)` is equivalent to `(x ** y) % z`. Here's an example with $x = 2$, $y = 3$, and $z = 2$:

```
>>> pow(2, 3, 2)
0
```

First, 2 is raised to the power 3 to get 8. Then $8 \% 2$ is calculated, which is 0 because 2 divides 8 with no remainder.

Check if a Float Is Integral

In Chapter 3 you learned about string methods like `.lower()`, `.upper()`, and `.find()`. Integers and floating-point numbers also have methods.

Number methods aren't used very often, but there is one that can be useful. Floating-point numbers have an `.is_integer()` method that returns `True` if the number is **integral**—meaning it has no fractional part—and returns `False` otherwise:

```
>>> num = 2.5
>>> num.is_integer()
False

>>> num = 2.0
>>> num.is_integer()
True
```

The `.is_integer()` method can be useful for validating user input. For example, if you are writing an app for a shopping cart for a store that sells pizzas, you will want to check that the quantity of pizzas the customer inputs is a whole number. You'll learn how to do these kinds of checks in Chapter 8.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Write a script that asks the user to input a number and then displays that number rounded to two decimal places. When run, your program should look like this:

```
Enter a number: 5.432
5.432 rounded to 2 decimal places is 5.43
```

2. Write a script that asks the user to input a number and then displays the absolute value of that number. When run, your program should look like this:

```
Enter a number: -10
The absolute value of -10 is 10.0
```

3. Write a script that asks the user to input two numbers by using the `input()` function twice, then display whether or not the difference between those two numbers is an integer. When run, your program

should look like this:

```
Enter a number: 1.5
```

```
Enter another number: .5
```

```
The difference between 1.5 and .5 is an integer? True!
```

If the user inputs two numbers whose difference is not integral, the output should look like this:

```
Enter a number: 1.5
```

```
Enter another number: 1.0
```

```
The difference between 1.5 and 1.0 is an integer? False!
```

[Leave feedback on this section »](#)

5.6 Print Numbers in Style

Displaying numbers to a user requires inserting numbers into a string. In Chapter 3, you learned how to do this with f-strings by surrounding a variable assigned to a number with curly braces:

```
>>> n = 7.125
>>> f"The value of n is {n}"
'The value of n is 7.125'
```

Those curly braces support a simple [formatting language](#) you can use to alter the appearance of the value in the final formatted string.

For example, to format the value of `n` in the above example to two decimal places, replace the contents of the curly braces in the f-string with `{n:.2f}`:

```
>>> n = 7.125
>>> f"The value of n is {n:.2f}"
'The value of n is 7.12'
```

The colon (:) after the variable `n` indicates that everything after it is part of the formatting specification. In this example, the formatting specification is `.2f`.

The `.2` in `.2f` rounds the number to two decimal places, and the `f` tells

Python to display `n` as a **fixed-point number**. This means the number is displayed with exactly two decimal places, even if the original number has fewer decimal places.

When `n = 7.125`, the result of `{n:.2f}` is `7.12`. Just like `round()`, Python rounds ties to even when formatting numbers inside of strings. So, if you replace `n = 7.125` with `n = 7.126`, then the result of `{n:.2f}` is `"7.13"`:

```
>>> n = 7.126
>>> f"The value of n is {n:.2f}"
'The value of n is 7.13'
```

To round to one decimal places, replace `.2` with `.1`:

```
>>> n = 7.126
>>> f"The value of n is {n:.1f}"
'The value of n is 7.1'
```

When you format a number as fixed-point, it's always displayed with the precise number of decimal places specified:

```
>>> n = 1
>>> f"The value of n is {n:.2f}"
'The value of n is 1.00'
>>> f"The value of n is {n:.3f}"
'The value of n is 1.000'
```

You can insert commas to group the integer part of large numbers by the thousands with the `,` option:

```
>>> n = 1234567890
>>> f"The value of n is {n:,}"
'The value of n is 1,234,567,890'
```

To round to some number of decimal places and also group by thousands, put the `,` before the dot `.` in your formatting specification:

```
>>> n = 1234.56
>>> f"The value of n is {n:.2f}"
'The value of n is 1,234.56'
```

The specifier `.2f` is useful for displaying currency values:

```
>>> balance = 2000.0
>>> spent = 256.35
>>> remaining = balance - spent

>>> f"After spending ${spent:.2f}, I was left with ${remaining:.2f}"
'After spending $256.35, I was left with $1,743.65'
```

Another useful option is `%`, which is used to display percentages. The `%` option multiplies a number by 100 and displays it in fixed-point format, followed by a percentage sign.

The `%` option should always go at the end of your formatting specification, and you can't mix it with the `f` option. For example, `.1%` displays a number as a percentage with exactly one decimal place:

```
>>> ratio = 0.9
>>> f"Over {ratio:.1%} of Pythonistas say 'Real Python rocks! ''"
'Over 90.0% of Pythonistas say 'Real Python rocks! '''

>>> # Display percentage with 2 decimal places
>>> f"Over {ratio:.2%} of Pythonistas say 'Real Python rocks! ''"
'Over 90.00% of Pythonistas say 'Real Python rocks! '''
```

The formatting mini language is powerful and extensive. You've only seen the basics here. For more information, you are encouraged to read the [official documentation](#).

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Print the result of the calculation `3 ** .125` as a fixed-point number with three decimal places.
2. Print the number `150000` as currency, with the thousands grouped with commas. Currency should be displayed with two decimal places.
3. Print the result of `2 / 10` as a percentage with no decimal places. The output should look like `20%`.

[Leave feedback on this section »](#)

5.7 Complex Numbers

Python is one of the few programming languages that provides built-in support for complex numbers. While complex numbers do not come up often outside the domains of scientific computing and computer graphics, Python's support for them is one of it's strengths.

Note

Complex numbers only come up in a few specific situations. Many programmers never need to use them.

Feel free to skip this section all together if you have no interest in how to work with complex numbers in Python. No other part of the book depends on the information in this section.

If you have ever taken a pre-calculus or higher level algebra math class, you may remember that a complex number is a number with two distinct components: a **real** component and an **imaginary** component.

There are several ways to denote complex numbers, but a common method is to indicate the real component with the letter *i* and the imaginary component with the letter *j*. For example, `1i + 2j` is the complex number with real part 1 and imaginary part 2.

To create a complex number in Python, you simply write the real part, followed by a plus sign and the imaginary part with the letter *j* at the

end:

```
>>> n = 1 + 2j
```

When you inspect the value of `n`, you'll notice that Python wraps the number with parentheses:

```
>>> n  
(1+2j)
```

This convention helps eliminate any confusion that the displayed output may represent a string or a mathematical expression.

Imaginary numbers come with two properties, `.real` and `.imag`, that return the real and imaginary component of the number, respectively:

```
>>> n.real  
1.0  
  
>>> n.imag  
2.0
```

Notice that Python returns both the real and imaginary components as floats, even though they were specified as integers.

Complex numbers also have a `.conjugate()` method that returns the complex conjugate of the number:

```
>>> n.conjugate()  
(1-2j)
```

For any complex number, its **conjugate** is the complex number with the same real part and an imaginary part that is the same in absolute value but with the opposite sign. So in this case, the complex conjugate of $1 + 2j$ is $1 - 2j$.

Note

The `.real` and `.imag` properties don't need parentheses after them like the method `.conjugate()` does.

The `.conjugate()` method is a function that performs an action on the complex number. `.real` and `.imag` don't perform any action, they just return some information about the number.

The distinction between methods and properties is a part of **object oriented programming**, which you will learn about in Chapter 10.

All of the arithmetic operators that work with floats and integers work with complex numbers also, except for the floor division (`//`) operator. Since this isn't a math book, we won't discuss the mechanics of **complex arithmetic**. Instead, here are some examples of using complex numbers with arithmetic operators:

```
>>> a = 1 + 2j
>>> b = 3 - 4j

>>> a + b
(4-2j)

>>> a - b
(-2+6j)

>>> a * b
(11+2j)

>>> a ** b
(932.1391946432212+95.9465336603419j)

>>> a / b
(-0.2+0.4j)
```

```
>>> a // b
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
TypeError: can't take floor of complex number.
```

Interestingly, although not surprising from a mathematical point of view, `int` and `float` objects also have the `.real` and `.imag` properties, as well as the `.conjugate()` method:

```
>>> x = 42
>>> x.real
42
>>> x.imag
0
>>> x.conjugate()
42

>>> y = 3.14
>>> y.real
3.14
>>> y.imag
0.0
>>> y.conjugate()
3.14
```

For floats and integers, `.real` and `.conjugate()` always return the number itself, and `.imag` always returns 0. One thing to notice, however, is that `n.real` and `n.imag` return an integer if `n` is an integer, and a float if `n` is a float.

Now that you have seen the basics of complex numbers, you might be wondering when you would ever need to use them. If you are learning Python for web development or automation, the truth is you may *never* need to use complex numbers.

On the other hand, complex numbers are important in domains such as scientific computing and computer graphics. If you ever work in those domains, you may find Python's built-in support for complex

numbers useful.

A detailed look at those topics is beyond the scope of this book. However, you will get an introduction to the NumPy package, a common tool for scientific computing with Python, in Chapter 17.

[Leave feedback on this section »](#)

5.8 Summary and Additional Resources

In this chapter you learned all about working with numbers in Python. You saw that there are two basic types of numbers—integers and floating-point numbers—and that Python also has built-in support for complex numbers.

First you learned how to do basic arithmetic with numbers using the `+`, `-`, `*`, `/`, and `%` operators. You saw how to write arithmetic expressions, and learned what the best practices are in [PEP 8](#) for formatting arithmetic expressions in your code.

Then you learned about floating-point numbers and how they may not always be 100% accurate. This limitation has nothing to do with Python. It is a fact of modern-day computing and is due to the way floating-point numbers are stored in a computer’s memory.

Next you saw how to round numbers to a given decimal place with the `round()` function, and learned that `round()` rounds ties to even, which is different from the way most people learned to round numbers in school.

Finally, you saw numerous ways to format numbers for display.

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-5

Additional Resources

To learn more, check out these resources:

- [Basic Data Types in Python](#)
- [How to Round Numbers in Python](#)
- Recommended resources on [realpython.com](#)

[Leave feedback on this section »](#)

Chapter 6

Functions and Loops

Functions are the building blocks of almost every Python program. They're where the real action takes place!

You've already seen how to use several functions, including `print()`, `len()`, and `round()`. These are all **built-in functions** because they come built into the Python language itself. You can also create **user-defined functions** that perform specific tasks.

Functions break code into smaller chunks, and are great for tasks that a program uses repeatedly. Instead of writing the same code each time the program need to perform the task, just call the function!

But sometimes you need to repeat some code several times in a row, and this is where **loops** come in.

In this chapter, you will learn:

- How to create user-defined functions
- How to write `for` and `while` loops
- What scope is and why it is important

Let's dive in!

[Leave feedback on this section »](#)

6.1 What is a Function, Really?

In the past few chapters you used functions like `print()` and `len()` to display text and determine the length of a string. But what is a function, really?

In this section you'll take a closer look at `len()` to learn more about what a function is and how it is executed.

Functions Are Values

One of the most important properties of a function in Python is that functions are values and can be assigned to a variable.

In IDLE's interactive window, inspect the name `len` by typing the following in at the prompt:

```
>>> len
<built-in function len>
```

When you hit `Enter`, Python tells you that the name `len` is a variable whose value is a built-in function.

Just like integer values have a type called `int`, and strings have a type `str`, function values also have a type:

```
>>> type(len)
<class 'builtin_function_or_method'>
```

Like any other variable, you can assign any value you want to `len`:

```
>>> len = "I'm not the len you're looking for."
>>> len
"I'm not the len you're looking for."
```

Now `len` has a string value, and you can verify that the type is `str` with `type()`:

```
>>> type(len)
<class 'str'>
```

The variable name `len` is a keyword in Python, and even though you can change its value, it's usually a bad idea to do so. Changing the value of `len` can make your code confusing because it's easy to mistake the new `len` for the built-in function.

Important

If you typed in the previous code examples, **you no longer have access to the built-in `len` function in IDLE.**

You can get it back with the following code:

```
>>> del len
```

The `del` keyword is used to un-assign a variable from a value. `del` stands for delete, but it doesn't delete the value. Instead, it detaches the name from the value and deletes the name.

Normally, after using `del`, trying to use the deleted variable name raises a `NameError`. In this case, however, the name `len` doesn't get deleted:

```
>>> len
<built-in function len>
```

Because `len` is a built-in function name, it gets reassigned to the original function value.

By going through each of these steps, we've seen that a function's name is separate from the function itself.

How Python Executes Functions

Now let's take a closer look at how Python executes a function.

The first thing to notice is that you can't execute a function by just

typing its name. You must **call** the function to tell Python to actually execute it.

Let's look at how this works with `len()`:

```
>>> # Typing just the name doesn't execute the function.  
>>> # IDLE inspects the variable as usual.  
>>> len  
<built-in function len>  
  
>>> # Use parentheses to call the function.  
>>> len()  
Traceback (most recent call last):  
  File "<pyshell#3>", line 1, in <module>  
    len()  
TypeError: len() takes exactly one argument (0 given)
```

In this example, Python raises a `TypeError` when `len()` is called because `len()` expects an argument.

An **argument** is a value that gets **passed** to the function as input. Some functions can be called with no arguments, and some can take as many arguments as you like. `len()` requires exactly one argument.

When a function is done executing, it **returns** a value as output. The return value usually — but not always — depends on the values of any arguments passed to the function.

The process for executing a function can be summarized in three steps:

1. The function is **called**, and any arguments are passed to the function as input.
2. The function **executes**, and some action is performed with the arguments.
3. The function **returns**, and the original function call is replaced with the return value.

Let's look at this in practice and see how Python executes the following line of code:

```
>>> num_letters = len("four")
```

First, `len()` is called with the argument "four". The length of the string "four" is calculated, which is the number 4. Then `len()` returns the number 4 and replaces the function call with the value.

So, after the function executes, the line of code looks like this:

```
>>> num_letters = 4
```

Then Python assigns the value 4 to `num_letters` and continues executing any remaining lines of code in the program.

Functions Can Have Side Effects

You've learned how to call a function and that they return a value when they are done executing. Sometimes, though, functions do more than just return a value.

When a function changes or affects something external to the function itself, it is said to have a **side effect**. You have already seen one function with a side effect: `print()`.

When you call `print()` with a string argument, the string is displayed in the Python shell as text. But `print()` doesn't return any text as a value.

To see what `print()` returns, you can assign the return value of `print()` to a variable:

```
>>> return_value = print("What do I return?")
What do I return?
>>> return_value
>>>
```

When you assign `print("What do I return?")` to `return_value`, the string

"What do I return?" is displayed. However, when you inspect the value of `return_value`, nothing is shown.

`print()` returns a special value called `None` that indicates the absence of data. `None` has a type called `NoneType`:

```
>>> type(return_value)
<class 'NoneType'>
>>> print(return_value)
None
```

When you call `print()`, the text that gets displayed is not the return value. It is a side effect of `print()`.

Now that you know that functions are values, just like strings and numbers, and have learned how functions are called and executed, let's take a look at how you can create your own user-defined functions.

[Leave feedback on this section »](#)

6.2 Write Your Own Functions

As you write longer and more complex programs, you may find that you need to use the same few lines of code repeatedly. Or maybe you need to calculate the same formula with different values several times in your code.

You might be tempted to copy and paste similar code to other parts of your program and modify it as needed, but this is usually a bad idea!

Repetitive code can be a nightmare to maintain. If you find a mistake in some code that's been copied and pasted all over the place, you'll end up having to apply the fix everywhere the code was copied. That's a lot of work, and you might miss a spot!

In this section, you'll learn how to define your own functions so that you can avoid repeating yourself when you need to reuse code. Let's

go!

The Anatomy of a Function

Every function has two parts:

1. The **function signature** defines the name of the function and any inputs it expects.
2. The **function body** contains the code that runs every time the function is used.

Let's start by writing a function that takes two numbers as input and returns their product. Here's what this function might look like, with the signature, body, and return statement identified with comments:

```
def multiply(x, y):  # Function signature
    # Function body
    product = x * y
    return product
```

It might seem odd to make a function for something as simple as the `*` operator. In fact, `multiply` is not a function you would probably write in a real-world scenario. But it makes a great first example for understanding how functions are created!

Let's break the function down to see what's going on.

The Function Signature

The first line of code in a function is called the **function signature**. It always starts with the `def` keyword, which is short for “define.”

Let's look more closely at the signature of the `multiply` function:

```
def multiply(x, y):
```

The function signature has four parts:

1. The `def` keyword

2. The function name, `multiply`
3. The parameter list, `(x, y)`
4. A colon (`:`) at the end of the line

When Python reads a line beginning with the `def` keyword, it creates a new function. The function is assigned to a variable with the same name as the function name.

Note

Since function names become variables, they must follow the same rules for variable names that you learned in Chapter 3.

So, a function name can only contain numbers, letters, and underscores, and must not begin with a number.

The parameter list is a list of parameter names surrounded by opening and closing parentheses. It defines the function's expected inputs. `(x, y)` is the parameter list for the `multiply` function. It creates two parameters, `x` and `y`.

A **parameter** is sort of like a variable, except that it has no value. It is a placeholder for actual values that are provided whenever the function is called with one or more arguments.

Code in the function body can use parameters as if they are variables with real values. For example, the function body may contain a line of code with the expression `x * y`.

Since `x` and `y` have no value, `x * y` has no value. Python saves the expression as a template and fills in the missing values when the function is executed.

A function can have any number of parameters, including no parameters at all!

The Function Body

The **function body** is the code that gets run whenever the function is used in your program. Here's the function body for the `multiply` function:

```
def multiply(x, y):
    # Function body
    product = x * y
    return product
```

`multiply` is a pretty simple function. Its body has only two lines of code!

The first line creates a variable called `product` and assigns to it the value `x * y`. Since `x` and `y` have no values yet, this line is really a template for the value `product` is assigned when the function is executed.

The second line of code is called a **return statement**. It starts with the `return` keyword and is followed by the variable `product`. When Python reaches the return statement, it stops running the function and returns the value of `product`.

Notice that both lines of code in the function body are indented. This is vitally important! Every line that is indented below the function signature is understood to be part of the function body.

For instance, the `print()` function in the following example is not a part of the function body because it is not indented:

```
def multiply(x, y):
    product = x * y
    return product

print("Where am I?") # Not in the function body.
```

If `print()` is indented, then it becomes a part of the function body even if there is a blank line between `print()` and the previous line:

```
def multiply(x, y):
    product = x * y
    return product

print("Where am I?") # In the function body.
```

There is one rule that you must follow when indenting code in a function's body. Every line must be indented by the same number of spaces.

Try saving the following code to a file called `multiply.py` and running it from IDLE:

```
def multiply(x, y):
    product = x * y
    return product # Indented with one extra space.
```

IDLE won't run the code! A dialog box appears with the error "unexpected indent." Python wasn't expecting the return statement to be indented differently than the line before it.

Another error occurs when a line of code is indented less than the line above it, but the indentation doesn't match any previous lines. Modify the `multiply.py` file to look like this:

```
def multiply(x, y):
    product = x * y
    return product # Indented less than previous line.
```

Now save and run the file. IDLE stops it with the error "unindent does not match any outer indentation level." The return statement isn't indented with the same number of spaces as any other line in the function body.

Note

Although Python has no rules for the number of spaces used to indent code in a function body, [PEP 8 recommends indenting with four spaces](#).

We follow this convention throughout this book.

Once Python executes a `return` statement, the function stops running and returns the value. If any code appears below the `return` statement that is indented so as to be part of the function body, it will never run.

For instance, the `print()` function will never be executed in the following function:

```
def multiply(x, y):
    product = x * y
    return product
    print("You can't see me!")
```

If you call this version of `multiply()`, you will never see the string "You can't see me!" displayed.

Calling a User-Defined Function

You call a user-defined function just like any other function. Type the function name followed by a list of arguments in between parentheses.

For instance, to call `multiply()` with the argument 2 and 4, just type:

```
multiply(2, 4)
```

Unlike built-in functions, user-defined functions are not available until they have been defined with the `def` keyword. You must define the function before you call it.

Try saving and running the following script:

```
num = multiply(2, 4)
print(num)

def multiply(x, y):
    product = x * y
    return product
```

When Python reads the line `num = multiply(2, 4)`, it doesn't recognize the name `multiply` and raises a `NameError`:

```
Traceback (most recent call last):
  File "C:Usersdaveamultiply.py", line 1, in <module>
    num = multiply(2, 4)
NameError: name 'multiply' is not defined
```

To fix the error, move the function definition to the top of the file:

```
def multiply(x, y):
    product = x * y
    return product

num = multiply(2, 4)
print(num)
```

Now when you save and run the script, the value `8` is displayed in the interactive window.

Functions With No Return Statement

All functions in Python return a value, even if that value is `None`. However, not all functions need a `return` statement.

For example, the following function is perfectly valid:

```
def greet(name):
    print(f"Hello, {name}!")
```

`greet()` has no `return` statement, but works just fine:

```
>>> greet("Dave")
Hello, Dave!
```

Even though `greet()` has no return statement, it still returns a value:

```
>>> return_value = greet("Dave")
Hello, Dave!
>>> print(return_value)
None
```

Notice also that the string "Hello, Dave!" is printed even when the result of `greet("Dave")` is assigned to a variable. That's because the call to `print()` inside of the `greet()` function body produces the side effect of always printing to the console.

If you weren't expecting to see "Hello, Dave!" printed, then you just experienced one of the issues with side effects. They aren't always expected!

When you create your own functions, you should always document what they do. That way other developers can read the documentation and know how to use the function and what to expect when it is called.

Documenting Your Functions

To get help with a function in IDLE's interactive window, you can use the `help()` function:

```
>>> help(len)
Help on built-in function len in module builtins:

len(obj, /)
    Return the number of items in a container.
```

When you pass a variable name or function name to `help()`, it displays some useful information about it. In this case, `help()` tells you that `len` is a built-in function that returns the number of items in a container.

Note

A **container** is a special name for an object that contains other objects. A string is a container because it contains characters.

You will learn about other container types in Chapter 9.

Let's see what happens when you call `help()` on the `multiply()` function:

```
>>> help(multiply)
Help on function multiply in module __main__:

multiply(x, y)
```

`help()` displays the function signature, but there isn't any information about what the function does. To better document `multiply()`, we need to provide a docstring.

A **docstring** is a triple-quoted string literal placed at the top of the function body. Docstrings are used to document what a function does and what kinds of parameters it expects.

Here's what `multiply()` looks like with a docstring added to it:

```
def multiply(x, y):
    """Return the product of two numbers x and y."""
    product = x * y
    return product
```

Update the `multiply.py` script with the docstring, then save and run the script. Now you can use `help()` in the interactive window to see the docstring:

```
>>> help(multiply)
Help on function multiply in module __main__:

multiply(x, y)
    Return the product of two numbers x and y.
```

PEP 8 doesn't say much about docstrings, except that [every function should have one](#).

There are a number of standardized docstring formats, but we won't get into them here. Some general guidelines for writing docstrings can be found in [PEP 257](#).

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Write a function called `cube()` with one number parameter and returns the value of that number raised to the third power. Test the function by displaying the result of calling your `cube()` function on a few different numbers.
2. Write a function called `greet()` that takes one string parameter called `name` and displays the text "Hello <name>!", where `<name>` is replaced with the value of the `name` parameter.

[Leave feedback on this section »](#)

6.3 Challenge: Convert Temperatures

Write a script called `temperature.py` that defines two functions:

1. `convert_cel_to_far()` which takes one `float` parameter representing degrees Celsius and returns a float representing the same temperature in degrees Fahrenheit using the following formula:

$$F = C * 9/5 + 32$$

2. `convert_far_to_cel()` which take one `float` parameter representing degrees Fahrenheit and returns a float representing the same temperature in degrees Celsius using the following formula:

$$C = (F - 32) * 5/9$$

The script should first prompt the user to enter a temperature in degrees Fahrenheit and then display the temperature converted to Celsius.

Then prompt the user to enter a temperature in degrees Celsius and display the temperature converted to Fahrenheit.

All converted temperatures should be rounded to 2 decimal places.

Here's a sample run of the program:

```
Enter a temperature in degrees F: 72  
72 degrees F = 22.22 degrees C
```

```
Enter a temperature in degrees C: 37  
37 degrees C = 98.60 degrees F
```

You can find the solutions to this code challenge and many other bonus resources online at realpython.com/python-basics/resources.

[Leave feedback on this section »](#)

6.4 Run in Circles

One of the great things about computers is that you can make them do the same thing over and over again, and they rarely complain or get tired.

A **loop** is a block of code that gets repeated over and over again either a specified number of times or until some condition is met. There are two kinds of loops in Python: **while loops** and **for loops**. In this section, you'll learn how to use both.

Let's start by looking at how `while` loops work.

The `while` Loop

`while` loops repeat a section of code while some condition is true. There are two parts to every `while` loop:

1. The **while statement** starts with the `while` keyword, followed by a **test condition**, and ends with a colon (:).
2. The **loop body** contains the code that gets repeated at each step of the loop. Each line is indented four spaces.

When a `while` loop is executed, Python evaluates the test condition and determines if it is true or false. If the test condition is true, then the code in the loop body is executed. Otherwise, the code in the body is skipped and the rest of the program is executed.

If the test condition is true and the body of the loop is executed, then once Python reaches the end of the body, it returns to the `while` statement and re-evaluates the test condition. If the test condition is still true, the body is executed again. If it is false, the body is skipped.

This process repeats over and over until the test condition fails, causing Python to loop over the code in the body of the `while` loop.

Let's look at an example. Type the following code into the interactive window:

```
>>> n = 1
>>> while n < 5:
...     print(n)
...     n = n + 1
...
1
2
3
4
```

First, the integer 1 is assigned to the variable `n`. Then a `while` loop is created with the test condition `n < 5`, which checks whether or not the value of `n` is less than 5.

If `n` is less than 5, the body of the loop is executed. There are two lines of code in the loop body. In the first line, the value of `n` is printed on the screen, and then `n` is **incremented** by 1 in the second line.

The loop execution takes place in five steps, described in the following table:

Step #	Value of n	Test Condition	What Happens
1	1	<code>1 < 5</code> (true)	1 printed; n incremented to 2
2	2	<code>2 < 5</code> (true)	2 printed; n incremented to 3
3	3	<code>3 < 5</code> (true)	3 printed; n incremented to 4
4	4	<code>4 < 5</code> (true)	4 printed; n incremented to 5
5	5	<code>5 < 5</code> (false)	Nothing printed; loop ends.

If you aren't careful, you can create an **infinite loop**. This happens when the test condition is always true. An infinite loop never terminates. The loop body keeps repeating forever.

Here's an example of an infinite loop:

```
>>> n = 1
>>> while n < 5:
...     print(n)
...
```

The only difference between this `while` loop and the previous one is that `n` is never incremented in the loop body. At each step of the loop, `n` is equal to 1. That means the test condition `n < 5` is always true, and the number 1 is printed over and over again forever.

Note

Infinite loops aren't inherently bad. Sometimes they are exactly the kind of loop you need.

For example, code that interacts with hardware may use an infinite loop to constantly check whether or not a button or switch has been activated.

If you run a program that enters an infinite loop, you can force Python to quit by pressing `Ctrl+C`. Python stops running the program and raises a `KeyboardInterrupt` error:

```
Traceback (most recent call last):
  File "<pyshell#8>", line 2, in <module>
    print(n)
KeyboardInterrupt
```

Let's look at an example of a `while` loop in practice. One use of a `while` loop is to check whether or not user input meets some condition and, if not, repeatedly ask the user for new input until valid input is received.

For instance, the following program continuously asks a user for a positive number until a positive number is entered:

```
num = float(input("Enter a positive number: "))

while num <= 0:
    print("That's not a positive number!")
    num = float(input("Enter a positive number: "))
```

First, the user is prompted to enter a positive number. The test condition `num <= 0` determines whether or not `num` is less than or equal to 0.

If `num` is positive, then the test condition fails. The body of the loop is skipped and the program ends.

Otherwise, if `num` is 0 or negative, the body of the loop executes. The

user is notified that their input was incorrect, and they are prompted again to enter a positive number.

`while` loops are perfect for repeating a section of code while some condition is met. They aren't well-suited, however, for repeating a section of code a specific number of times.

The `for` Loop

A `for` loop executes a section of code once for each item in a collection of items. The number of times that the code is executed is determined by the number of items in the collection.

Like its `while` counterpart, the `for` loop has two main parts:

1. The **for statement** begins with the `for` keyword, followed by a **membership expression**, and ends in a colon (:).
2. The **loop body** contains the code to be executed at each step of the loop, and is indented four spaces.

Let's look at an example. The following `for` loop prints each letter of the string "Python" one at a time:

```
for letter in "Python":  
    print(letter)
```

In this example, the `for` statement is `for letter in "Python"`. The membership expression is `letter in "Python"`.

At each step of the loop, the variable `letter` is assigned the next letter in the string "Python", and then the value of `letter` is printed.

The loops runs once for each character in the string "Python", so the loop body executes six times. The following table summarizes the execution of this `for` loop:

Step #	Value of <code>letter</code>	What Happens
1	"P"	P is printed

Step #	Value of letter	What Happens
2	"y"	y is printed
3	"t"	t is printed
4	"h"	h is printed
5	"o"	o is printed
6	"n"	n is printed

To see why `for` loops are better for looping over collections of items, let's re-write the `for` loop in previous example as a `while` loop.

To do so, we can use a variable to store the index of the next character in the string. At each step of the loop, we'll print out the character at the current index and then increment the index.

The loop will stop once the value of the index variable is equal to the length of the string. Remember, indices start at 0, so the last index of the string "Python" is 5.

Here's how you might write that code:

```
word = "Python"
index = 0

while index < len(word):
    print(word[index])
    index = index + 1
```

That's significantly more complex than the `for` loop version!

Not only is the `for` loop less complex, the code itself looks more natural. It more closely resembles how you might describe the loop in English.

Note

You may sometimes hear people describe some code as being particularly “Pythonic.” The term **Pythonic** is generally used to describe code that is clear, concise, and uses Python’s built-in features to its advantage.

In these terms, using a `for` loop to loop over a collection of items is more Pythonic than using a `while` loop.

Sometimes it’s useful to loop over a range of numbers. Python has a handy built-in function `range()` that produces just that — a range of numbers!

For example, `range(3)` returns the range of integers starting with 0 and up to, but not including, 3. That is, `range(3)` is the range of numbers 0, 1, and 2.

You can use `range(n)`, where `n` is any positive number, to execute a loop exactly `n` times. For instance, the following `for` loop prints the string "Python" three times:

```
for n in range(3):
    print("Python")
```

You can also give a range a starting point. For example, `range(1, 5)` is the range of numbers 1, 2, 3, and 4. The first argument is the starting number, and the second argument is the endpoint, which is not included in the range.

Using the two-argument version of `range()`, the following `for` loop prints the square of every number starting with 10 and up to, but not including, 20:

```
for n in range(10, 20):
    print(n * n)
```

Let’s look at a practical example. The following program asks the user to input an amount and then displays how to split that amount be-

tween 2, 3, 4, and 5 people:

```
amount = float(input("Enter an amount: "))

for num_people in range(2, 6):
    print(f"{num_people} people: ${amount / num_people:.2f} each")
```

The `for` loop loops over the number 2, 3, 4, and 5, and prints the number of people and the amount each person should pay. The formatting specifier `.2f` is used to format the amount as fixed-point number rounded to two decimal places and commas every three digits.

Running the program with the input 10 produces the following output:

```
Enter an amount: 10
2 people: $5.00 each
3 people: $3.33 each
4 people: $2.50 each
5 people: $2.00 each
```

`for` loops are generally used more often than `while` loops in Python. Most of the time, a `for` loop is more concise and easier to read than an equivalent `while` loop.

Nested Loops

As long as you indent the code correctly, you can even put loops inside of other loops.

Type the following into IDLE's interactive window:

```
for n in range(1, 4):
    for j in range(4, 7):
        print(f"n = {n} and j = {j}")
```

When Python enters the body of the first `for` loop, the variable `n` is assigned the value 1. Then the body of the second `for` loop is executed and `j` is assigned the value 4. The first thing printed is `n = 1 and j = 4`.

After executing the `print()` function, Python returns to the *inner* `for` loop, assigns to `j` the value of 5, and then prints `n = 1` and `j = 5`. Python doesn't return the outer `for` loop because the inner `for` loop, which is inside the body of the outer `for` loop, isn't done executing.

Next, `j` is assigned the value 6 and Python prints `n = 1` and `j = 6`. At this point, the inner `for` loop is done executing, so control returns to the outer `for` loop.

The variable `n` gets assigned the value 2, and the inner `for` loop executes a second time. That is, `j` is assigned the value 4 and `n = 2` and `j = 4` is printed to the console.

The two loops continue to execute in this fashion, and the final output looks like this:

```
n = 1 and j = 4
n = 1 and j = 5
n = 1 and j = 6
n = 2 and j = 4
n = 2 and j = 5
n = 2 and j = 6
n = 3 and j = 4
n = 3 and j = 5
n = 3 and j = 6
```

A loop inside of another loop is called a **nested loop**, and they come up more often than you might expect. You can nest `while` loops inside of `for` loops, and vice versa, and even nest loops more than two levels deep!

Important

Nesting loops inherently increases the complexity of your code, as you can see by the dramatic increase in the number of steps run in the previous example compared to examples with a single `for` loop.

Using nested loops is sometimes the only way to get something done, but too many nested loops can have a negative effect on a program's performance.

Loops are a powerful tool. They tap into one of the greatest advantages computers provide as tools for computation: the ability to repeat the same task a vast number of times without tiring and without complaining.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Write a `for` loop that prints out the integers 2 through 10, each on a new line, by using the `range()` function.
2. Use a `while` loop that prints out the integers 2 through 10 (Hint: You'll need to create a new integer first.)
3. Write a function called `doubles()` that takes one number as its input and doubles that number. Then use the `doubles()` function in a loop to double the number 2 three times, displaying each result on a separate line. Here is some sample output:

```
4  
8  
16
```

[Leave feedback on this section »](#)

6.5 Challenge: Track Your Investments

In this challenge, you will write a program called `invest.py` that tracks the growing amount of an investment over time.

An initial deposit, called the principal amount, is made. Each year, the amount increases by a fixed percentage, called the annual rate of return.

For example, a principal amount of \$100 with an annual rate of return of 5% increases the first year by \$5. The second year, the increase is 5% of the new amount \$105, which is \$5.25.

Write a function called `invest` with three parameters: the principal amount, the annual rate of return, and the number of years to calculate. The function signature might look something like this:

```
def invest(amount, rate, years):
```

The function then prints out the amount of the investment, rounded to 2 decimal places, at the end of each year for the specified number of years.

For example, calling `invest(100, .05, 4)` should print the following:

```
year 1: $105.00
year 2: $110.25
year 3: $115.76
year 4: $121.55
```

To finish the program, prompt the user to enter an initial amount, an annual percentage rate, and a number of years. Then call `invest()` to display the calculations for the values entered by the user.

You can find the solutions to this code challenge and many other bonus resources online at realpython.com/python-basics/resources.

[Leave feedback on this section »](#)

6.6 Understand Scope in Python

Any discussion of functions and loops in Python would be incomplete without some mention of the issue of **scope**.

Scope can be one of the more difficult concepts to understand in programming, so in this section you will get a gentle introduction to it.

By the end of this section, you'll know what a **scope** is and why it is important. You will also learn the LEGB rule for **scope resolution**.

What Is a Scope?

When you assign a value to a variable, you are giving that value a name. Names are unique. For example, you can't assign the same name to two different numbers.

```
>>> x = 2
>>> x
2

>>> x = 3
>>> x
3
```

When you assign 3 to x, you can no longer recall the value 2 with the name x.

This behavior makes sense. After all, if the variable x has the values 2 and 3 simultaneously, how do you evaluate $x + 2$? Should it be 4 or 5?

As it turns out, there *is* a way to assign the same name to two different values. Try running the following script:

```
x = "Hello World"

def func():
    x = 2
```

```
print(f"Inside 'func', x has the value {x}")

func()
print(f"Outside 'func', x has the value {x}")
```

In this example, the variable `x` is assigned two different values. `x` is assigned "Hello, World" at the beginning, and is assigned 2 inside of `func()`.

The output of the script, which you might find surprising, looks like this:

```
Inside 'func', x has the value 2
Outside 'func', x has the value Hello World
```

How does `x` still have the value "Hello World" after calling `func()`, which changes the value of `x` to 2?

The answer is that the function `func()` has a different **scope** than the code that exists outside of the function. That is, you can name an object inside `func()` the same name as something outside `func()` and Python can keep the two separated.

The function body has what is known as a **local scope**, with its own set of names available to it. Code outside of the function body is in the **global scope**.

You can think of a scope as a set of names mapped to objects. When you use a particular name in your code, such as a variable or a function name, Python checks the current scope to determine whether or not that name exists.

Scope Resolution

Scopes have a hierarchy. For example, consider the following:

```
x = 5
```

```
def outer_func():
    y = 3

    def inner_func():
        z = x + y
        return z

    return inner_func()
```

Note

The `inner_func()` function is called an **inner function** because it is defined inside of another function. Just like you can nest loops, you can also define functions within other functions!

You can read more about inner functions in Real Python's article [Inner Functions—What Are They Good For?](#).

The variable `z` is in the local scope of `inner_func()`. When Python executes the line `z = x + y`, it looks for the variables `x` and `y` in the local scope. However, neither of them exist there, so it moves up to the scope of the `outer_func()` function.

The scope for `outer_func()` is an **enclosing** scope of `inner_func()`. It is not quite the global scope, and is not the local scope for `inner_func()`. It lies in between those two.

The variable `y` is defined in the scope for `outer_func()` and is assigned the value 3. However, `x` does not exist in this scope, so Python moves up once again to the global scope. There it finds the name `x`, which has the value 5. Now that the names `x` and `y` are resolved, Python can execute the line `z = x + y`, which assigns to `z` the value of 8.

The LEGB Rule

A useful way to remember how Python resolves scope is with the **LEGB** rule. This rule is an acronym for **L**ocal, **E**nclosing, **G**lobal, **B**uilt-in.

Python resolves scope in the order in which each scope appears in the list **LEGB**. Here is a quick overview to help you remember how all of this works:

Local (L): The local, or current, scope. This could be the body of a function or the top-level scope of a script. It always represents the scope that the Python interpreter is currently working in.

Enclosing (E): The enclosing scope. This is the scope one level up from the local scope. If the local scope is an inner function, the enclosing scope is the scope of the outer function. If the scope is a top-level function, the enclosing scope is the same as the global scope.

Global (G): The global scope, which is the top-most scope in the script. This contains all of the names defined in the script that are not contained in a function body.

Built-in (B): The built-in scope contains all of the names, such as keywords, that are built-in to Python. Functions such as `round()` and `abs()` are in the built-in scope. Anything that you can use without first defining yourself is contained in the built-in scope.

Break the Rules

Consider the following script. What do you think the output is?

```
total = 0

def add_to_total(n):
    total = total + n

add_to_total(5)
print(total)
```

You would think that script outputs the value 5, right? Try running it to see what happens.

Something unexpected occurs. You get an error!

```
Traceback (most recent call last):
  File "C:/Users/davea/stuff/python/scope.py", line 6, in <module>
    add_to_total(5)
  File "C:/Users/davea/stuff/python/scope.py", line 4, in add_to_total
    total = total + n
UnboundLocalError: local variable 'total' referenced before assignment
```

Wait a minute! According to the LEGB rule, Python should have recognized that the name `total` doesn't exist in the `add_to_total()` function's local scope and moved up to the global scope to resolve the name, right?

The problem here is that the script attempts to make an assignment to the variable `total`, which creates a new name in the local scope. Then, when Python executes the right-hand side of the assignment it finds the name `total` in the local scope with nothing assigned to it yet.

These kinds of errors are tricky and are one of the reasons it is best to use unique variable and function names no matter what scope you are in.

You can get around this issue with the `global` keyword. Try running the following altered script:

```
total = 0

def add_to_total(n):
    global total
    total = total + n

add_to_total(5)
print(total)
```

This time, you get the expected output 5. Why's that?

The line `global total` tells Python to look in the global scope for the name `total`. That way, the line `total = total + n` does not create a new local variable.

Although this “fixes” the script, the use of the `global` keyword is considered bad form in general.

If you find yourself using `global` to fix problems like the one above, stop and think if there is a better way to write your code. Often, you’ll find that there is!

[Leave feedback on this section »](#)

6.7 Summary and Additional Resources

In this chapter, you learned about two of the most essential concepts in programming: functions and loops.

First, you learned how to define your own custom functions. You saw that functions are made up of two parts:

1. The **function signature**, which starts with the `def` keyword and includes the name of the function and the function’s parameters
2. The **function body**, which contains the code that runs whenever the function is called.

Functions help avoid repeating similar code throughout a program by creating re-usable components. This helps make code easier to read and maintain.

Then you learned about Python’s two kinds of loops:

1. **while loops** repeat some code while some condition remains true
2. **for loops** repeat some code for each element in a set of objects

Finally, you learned what a **scope** is and how Python resolves scope using the LEGB rule.

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-6

Additional Resources

To learn more about functions and loops, check out the following resources:

- Python “while” Loops (Indefinite Iteration)
- Python “for” Loops (Definite Iteration)
- Recommended resources on realpython.com

[Leave feedback on this section »](#)

Chapter 7

Finding and Fixing Code Bugs

Everyone makes mistakes—even seasoned professional developers!

IDLE is pretty good at catching mistakes like syntax and run-time errors, but there's a third type of error that you may have already experienced. **Logic errors** occur when an otherwise valid program doesn't do what was intended.

Logic errors cause unexpected behaviors called **bugs**. Removing bugs is called **debugging**, and a **debugger** is a tool that helps you hunt down bugs and understand why they are happening.

Knowing how to find and fix bugs in your code is a skill that you will use for your entire coding career!

In this chapter, you will:

- Learn how to use IDLE's Debug Control Window
- Practice debugging on a buggy function

Let's go!

[Leave feedback on this section »](#)

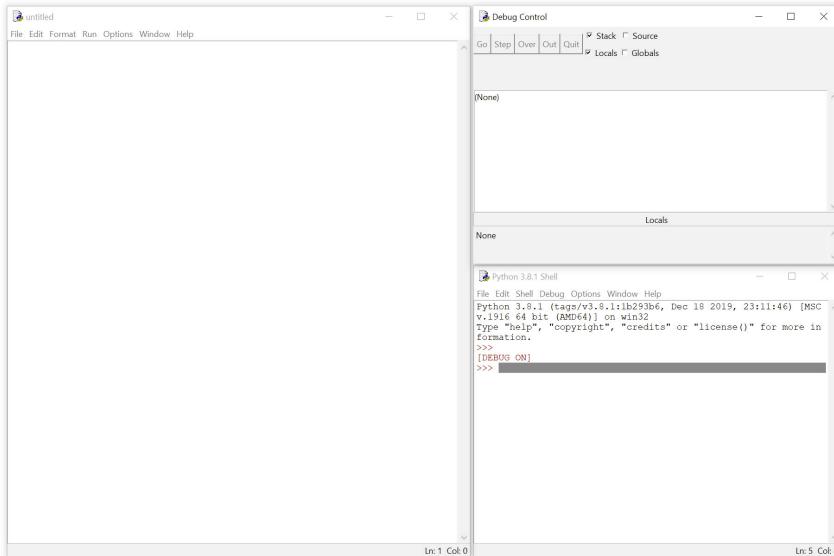
7.1 Use the Debug Control Window

The main interface to IDLE's debugger is through the Debug Control Window, which we'll refer to as the Debug window for short. You can open the Debug window by selecting `Debug > Debugger` from the menu in the interactive window. Go ahead and open the Debug window.

Note

If the `Debug` menu is missing from your menu bar, make sure the interactive window is in focus by clicking that window.

Open a new script window and arrange the three windows on your screen so that you can see all of them simultaneously. Here's one way you could rearrange the windows:



Note

Whenever the Debug window is open, the prompt in the interactive window has [DEBUG ON] next to it to indicate that the debugger is open.

In this section you'll learn how the Debug window is organized, how to step through your code with the debugger one line at a time, and how to set breakpoints to help speed up the debugging process.

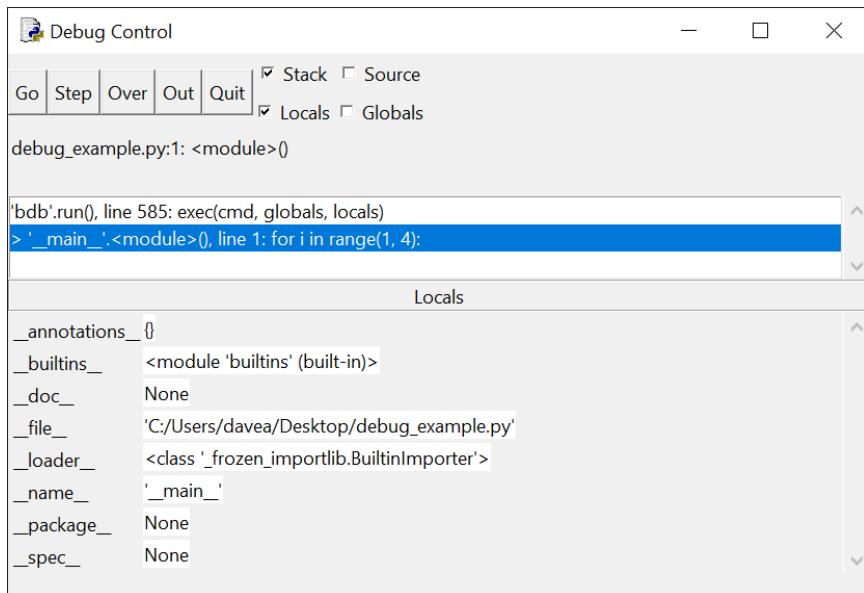
The Debug Control Window: An Overview

To see how the debugger works, let's start by writing a simple program without any bugs. Type the following into the script window:

```
for i in range(1, 4):
    j = i * 2
    print(f"i is {i} and j is {j}")
```

When you save and run this script with the Debug window open, you'll notice that execution doesn't get very far. The Debug Control window will look like this:

7.1. Use the Debug Control Window



Notice that the *Stack* panel at the top of the window contains the following message:

```
> '__main__'.<module>(), line 1: for i in range(1, 4):
```

This tells you that line 1 (which contains the code for `i in range(1, 4):`) is *about* to be run but has not started yet. The `'__main__'.module()` part of the message in the debugger refers to the fact that we're currently in the “main” section of the script, as opposed to being, for example, in a function definition before the main block of code has been reached.

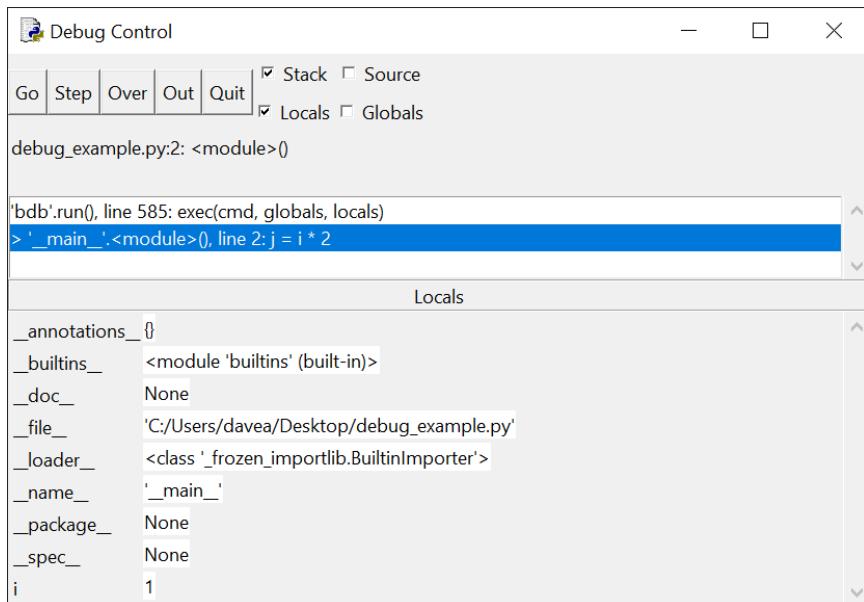
Below the *Stack* panel, there is a *Locals* panel that lists some strange looking stuff like `__annotations__`, `__builtins__`, `__doc__`, and so on. These are some internal system variables that you can ignore for now. As your program runs, you will see variables declared in the code displayed in this window so that you can keep track of their value.

There are five buttons located at the top left-hand corner of the Debug window: `Go`, `Step`, `Over`, `Out`, and `Quit`. These buttons control how the debugger moves through your code.

In the following sections, we'll explore what each of these buttons does, starting with the **Step** button.

The **Step** Button

Go ahead and click the **Step** button at the top left-hand corner of the Debug window. The Debug window changes a bit to look like this:



There are two differences to pay attention to here. First, the message in the *Stack* window changes to:

```
> '__main__.<module>(), line 2: j = i * 2:'
```

At this point, line 1 of your code was run, and the debugger has stopped just before executing line 2.

The second change to notice is the new variable *i* that is assigned the value 1 in the *Locals* panel. That's because the `for` loop in the first line of code created the variable *i* and assigned it the value 1.

Continue hitting the `Step` button to walk through your code line by line, watching what happens in the debugger window. When you arrive at the line `print(f'i is {i} and j is {j}')`, you can see the output displayed in the interactive window one piece at a time.

More importantly, you can track the growing values of `i` and `j` as you step through the `for` loop. You can probably imagine how beneficial this feature is when trying to locate the source of bugs in your programs. Knowing each variables value at each line of code can help you pinpoint where things go wrong.

Breakpoints and the “Go” Button

Often, you may know that the bug must be in a particular section of your code, but you may not know precisely where. Rather than clicking the `Step` button all day long, you can set a **breakpoint** that tells the debugger to run all code before the breakpoint continuously until the breakpoint is reached.

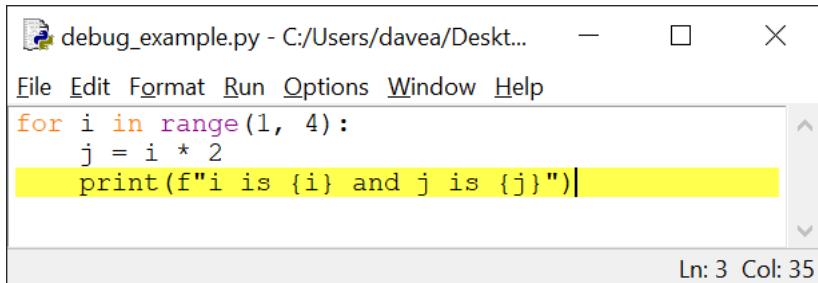
Breakpoints tell the debugger when to pause code execution so that you can take a look at the current state of the program. They don’t actually break anything.

To set a breakpoint, right-click (Mac: `Ctrl` + `Click`) on the line of code in your script window you would like to pause at and select `Set Breakpoint`. IDLE highlights the line in yellow to indicate that your breakpoint has been set. You can remove a breakpoint at any time by right-clicking on the line with a breakpoint and selecting `Clear Breakpoint`.

Go ahead and press the `Quit` button at the top of the Debug Control Window to turn off the debugger for now. This won’t close the window, and you’ll want to keep it open because you’ll be using it again in just a moment.

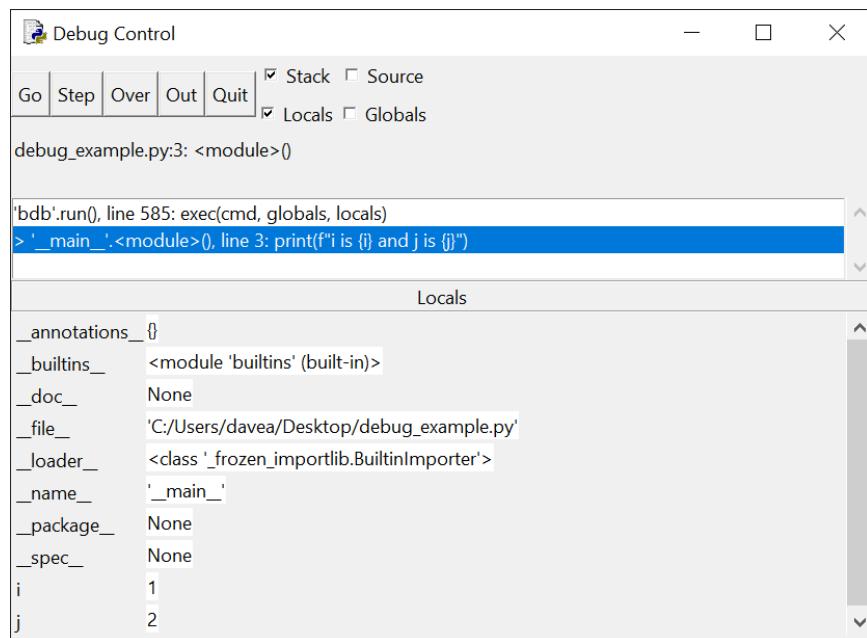
Set a breakpoint on the line of code with the `print()` statement. The script window should now look like this:

7.1. Use the Debug Control Window



```
debug_example.py - C:/Users/davea/Desktop/Debug Example  
File Edit Format Run Options Window Help  
for i in range(1, 4):  
    j = i * 2  
    print(f"i is {i} and j is {j}")|  
Ln: 3 Col: 35
```

Now run the script by pressing F5. Just like before, the *Stack* panel of the Debug Control Window indicates that debugger has started and is waiting to execute line 1. This time, instead of clicking on the **Step** button, click on the **Go** button and watch what happens to the Debug window:



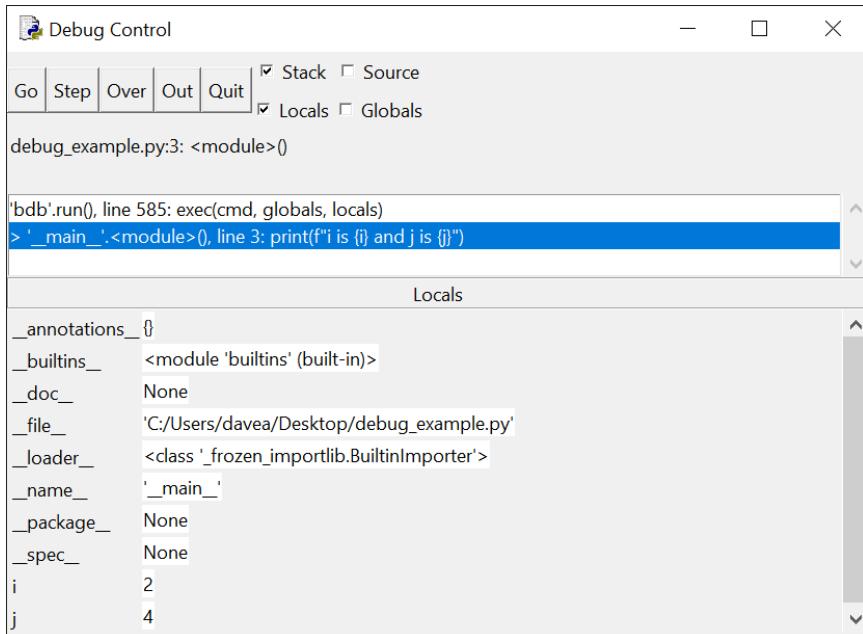
Locals	
annotations	{} None
builtins	<module 'builtins' (built-in)>
doc	None
file	'C:/Users/davea/Desktop/debug_example.py'
loader	<class '_frozen_importlib.BuiltinImporter'>
name	'_main_'
package	None
spec	None
i	1
j	2

The *Stack* panel now shows following message indicating that it is waiting to execute line 3:

7.1. Use the Debug Control Window

```
> '__main__'.<module>(), line 3: print(f'i is {i} and j is {j}')"
```

If you look at the *Locals* panel, you will see that both variable *i* and *j* have the values 1 and 2, respectively. By clicking on **Go**, you told the debugger to run your code continuously until reaching either a breakpoint or the end of the program. Now press “Go” again. The Debug window now looks like this:



Do you see what changed? The same message as before is displayed in the *Stack* panel, indicating the debugger is waiting to execute line 3 again. However, now the values of the variables *i* and *j* are 2 and 4. The interactive window also displays the output from running the line with `print()` in it the first time.

Each time you press the **Go** button, the debugger runs the code continuously until the next breakpoint is reached. Since you set the breakpoint on line 3, which is inside of the `for` loop, the debugger stops on this line each time it goes through the loop.

Press `Go` a third time. Now `i` and `j` have the values 3 and 6. What do you think happens when you press `Go` one more time? Since the `for` loop only iterates 3 times, when you press `Go` this time, the program finishes running.

“Over” and “Out”

The `Over` button works as sort of a combination of `Step` and `Go`. It steps over a function or loop. In other words, if you’re about to `Step` into a function with the debugger, you can still run that function’s code without having to `Step` all the way through each line of it. The `Over` button takes you directly to the result of running that function.

Likewise, if you’re already inside of a function or loop, the `Out` button executes the remaining code inside the function or loop body and then pauses.

In the next section, you’ll look at some buggy code and learn how to fix it with IDLE.

[Leave feedback on this section »](#)

7.2 Squash Some Bugs

Now that you’ve gotten comfortable using the Debug Control Window let’s take a look at a buggy program.

The following code defines a function `add_underscores()` that takes a single string object `word` as an argument and returns a new string containing a copy `word` with each character surrounded by underscores. For example, `add_underscores("python")` should return `"_p_y_t_h_o_n_"`.

Here’s the buggy code:

```
def add_underscores(word):
    new_word = "_"
    for i in range(0, len(word)):
```

```
new_word = word[i] + "_"
return new_word

phrase = "hello"
print(add_underscores(phrase))
```

Save and run the above script. The expected output is `_h_e_l_l_o_`, but instead all you see is `o_`, the letter "o" followed by a single underscore. If you already see what the problem with the code is, don't just fix it. The point of this section is to learn how to use IDLE's debugger to identify the problem. If you don't see what the problem is, don't worry! By the end of this section, you'll have found it and will be able to identify problems like it in other code you encounter.

Note

When working with real-world problems, debugging can often be difficult and time-consuming, and bugs can be subtle and hard to identify. While this section looks at a relatively simple bug, the important thing to take away from this is the methodology used to inspect the code.

Debugging is problem-solving, and as you become more experienced, you will develop your own approaches. In this section, you'll learn a simple four-step method to help get you started:

1. Guess which section of code may contain the bug.
2. Set a breakpoint and inspect the code by stepping through the buggy section one line at a time, keeping track of important variables along the way.
3. Identify the line of code, if any, with the error and make a change to solve the problem.
4. Repeat steps 1–3 as needed until the code works as expected.

Step 1: Make a Guess About Where the Bug Is Located

The first step is to identify the section of code that likely contains the bug. You may not be able to identify exactly where the bug is at first, but you can usually make a reasonable guess about which section of your code has an error.

Notice that the script is split into two distinct sections: a function definition (where `add_underscores()` is defined), and a “main” code block that defines a variable `phrase` with the value `"hello"` and then prints the result of calling `add_underscores(phrase)`.

Look at the “main” section:

```
phrase = "hello"  
print(add_underscores(phrase))
```

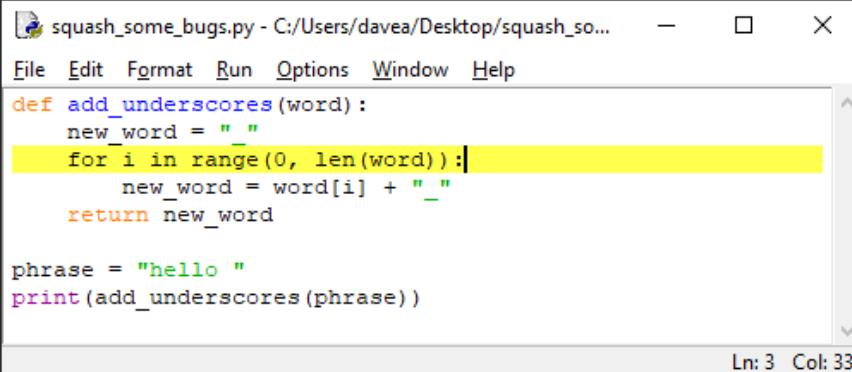
Do you think the problem could be here? It doesn’t look like it, right? Everything about those two lines of code looks good. So, the problem must be in the function definition:

```
def add_underscores(word):  
    new_word = "_"  
    for i in range(0, len(word)):  
        new_word = word[i] + "_"  
    return new_word
```

The first line of code inside the function creates a variable `new_word` with the value `_`. All good there, so we can conclude that the problem is somewhere in the body of the `for` loop.

Step 2: Set a Breakpoint and Inspect the Code

Now that you’ve identified where the bug must be, set a breakpoint at the start of the `for` loop so that you can trace out exactly what’s happening inside with the Debug Control Window:

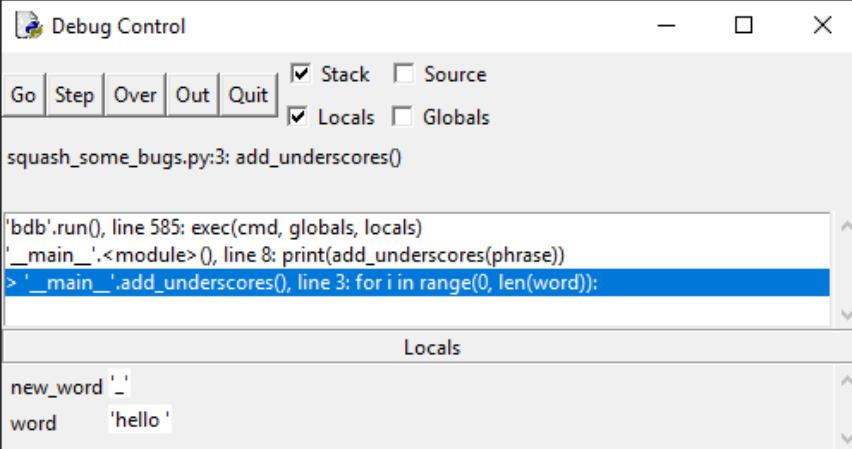


```
squash_some_bugs.py - C:/Users/davea/Desktop/squash_so...
File Edit Format Run Options Window Help
def add_underscores(word):
    new_word = "_"
    for i in range(0, len(word)):
        new_word = word[i] + "_"
    return new_word

phrase = "hello "
print(add_underscores(phrase))

Ln: 3 Col: 33
```

Now open the Debug Control Window and run the script. Execution still pauses on the very first line it sees (which is defining the function). Press the “Go” button to run through the code until the breakpoint is encountered. The Debug window will now look like this:



Debug Control

Stack Source
Locals Globals

squash_some_bugs.py:3: add_underscores()

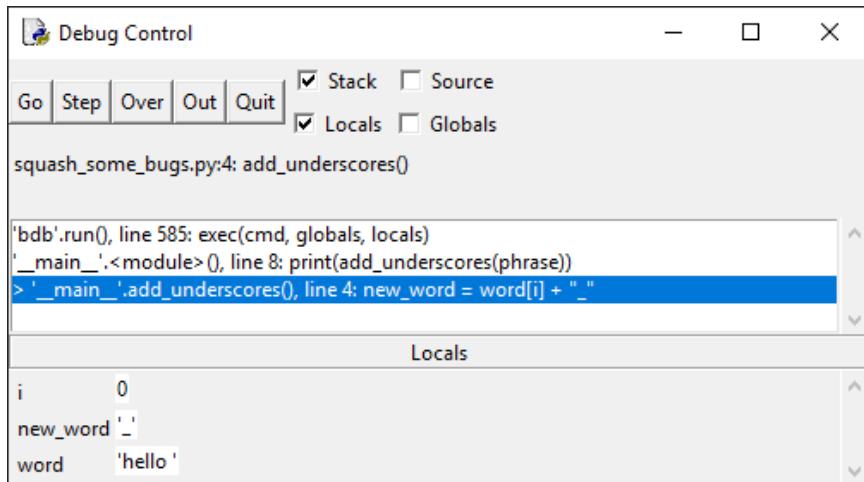
'bdb'.run(), line 585: exec(cmd, globals, locals)
'_main_'.<module>(), line 8: print(add_underscores(phrase))
> '_main_'.add_underscores(), line 3: for i in range(0, len(word)):

Locals

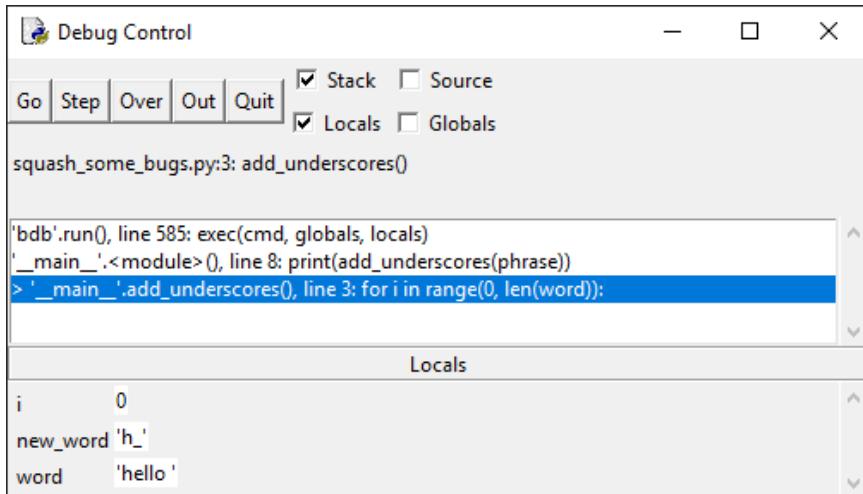
new_word '_'
word 'hello'

At this point, the code is paused just before entering the `for` loop in the `add_underscores()` function. Notice that two local variables, `word` and `new_word` are displayed in the `Locals` panel. Currently, `word` has the value "hello" and `new_word` has the value "_", as expected.

Click the **Step** button once to enter the `for` loop. The Debug window changes and a new variable `i` with the value `0` is displayed in the “Locals” panel. `i` is the counter used in the `for` loop, and you can use it to keep track of which iteration of the `for` loop you are currently looking at:



Click **Step** one more time. If you look at the *Locals* panel, you’ll see that the variable `new_word` has taken on the value "h_":



This isn't right. Originally, `new_word` had the value `"_"` and on the second iteration of the `for` loop it should now have the value `"_h_"`. If you click `Step` a few more times, you'll see that `new_word` gets set to `e_`, then `l_`, and so on.

Step 3: Identify the Error and Attempt to Fix It

The conclusion you can make at this point is that `new_word` is overwritten at each iteration of the `for` loop with the next character in the string `"hello"` and a trailing underscore. Since there is only one line of code inside the `for` loop, you know that the problem must be with the following code:

```
new_word = word[i] + "_"
```

Look at that closely. This line tells Python to get the next character of `word`, tack an underscore to the end of it, and assign this new string to the variable `new_word`. This is exactly the behavior you've witnessed by stepping through the `for` loop!

To fix the problem, you need to tell Python to concatenate the string `word[i] + "_"` to the existing value of `new_word`. Press the "Quit" button

in the Debug Control Window, but don't close that window just yet. Open the script window and change the line inside the `for` loop to:

```
new_word = new_word + word[i] + "_"
```

Step 4: Repeat Steps 1–3 Until the Bug is Gone

Save the new changes to your script and run it again. In the Debug window, press the `Go` button to execute the code up until the breakpoint.

Note

If you closed the debugger in the previous step without clicking on `Quit`, you may see the following error when re-opening the Debug Control Window:

```
You can only toggle the debugger when idle
```

Always be sure to click `Go` or `Quit` when you're finished with a debugging session instead of just closing the debugger, or you might have trouble reopening it. To get rid of this error, you'll have to close IDLE and re-open it.

Just like before, your script is now paused just before entering the `for` loop in the `add_underscores()` function. Press the `Step` button repeatedly and watch what happens to the `new_word` variable at each iteration. What do you see now? Success! Everything is working as expected!

In this example, your first attempt at fixing the bug worked, so you don't need to repeat steps 1–3 anymore. However, this won't always be the case. Sometimes you'll have to repeat the process several times before you've fixed a bug.

It's also important to keep in mind that tools like debuggers don't tell you how to fix a bug. They only help you identify where exactly a problem occurs in your code.

Alternative Ways to Find Mistakes in Your Code

Debugging can be tricky and time-consuming, but sometimes it's the most reliable way to find errors that you've overlooked. However, before you open a debugger, it is sometimes simpler to locate errors using well placed `print()` functions to display the values of your variables.

For example, instead of debugging the previous script with the Debug Control Window, you could add the following line to the end of the `for` loop in the `add_underscores()` function:

```
print(f"i = {i}; new_word = {new_word}")
```

The altered script would then look like this:

```
def add_underscores(word):
    new_word = "_"
    for i in range(0, len(word)):
        new_word = word[i] + "_"
        print(f"i = {i}; new_word = {new_word}")
    return new_word

phrase = "hello"
print(add_underscores(phrase))
```

When you run the script, the interactive window displays the following output:

```
i = 0; new_word = h_
i = 1; new_word = e_
i = 2; new_word = l_
i = 3; new_word = l_
i = 4; new_word = o_
o_
```

This shows you what the value of `new_word` is at each iteration of the `for` loop. The final line containing just a single underscore is the result of running `print(add_underscores(phrase))` at the end of the script.

By looking at the above output, you could come to the same conclusion you did while debugging with the Debug Control Window. That is, the problem is that `new_word` is overwritten at each iteration.

Many Python programmers prefer this simple method for some quick and dirty debugging on the fly. It is a handy technique but has some disadvantages when compared to IDLE’s debugger.

The most significant disadvantage is that debugging with the `print()` function requires you to run your entire script each time you want to inspect the values of your variables. For long scripts, this can be an enormous waste of time compared to setting breakpoints and using the “Go” button in the Debug Control Window.

Another disadvantage is that you’ll have to remember to remove those `print()` function calls from your code when you are done debugging it. Otherwise, users may see unnecessary and potentially confusing output when they run your program.

The example loop in this section may be a good example for illustrating the process of debugging, but it is not the best example of Pythonic code. The use of the index `i` is a giveaway that there might be a better way to write the loop.

One way to improve this loop is to iterate over the characters in the string `word` directly. Here’s one way to do that:

```
def add_underscores(word):
    new_word = "_"_
    for char in word:
        new_word = new_word + char + "_"
    return new_word
```

The process of re-writing existing code to be cleaner, easier to read and understand, or adhere to code standards set by a team is called **refactoring**. We won’t discuss refactoring much in this course, but it is an essential part of writing professional quality code.

[Leave feedback on this section »](#)

7.3 Summary and Additional Resources

In this chapter, you learned about IDLE’s Debug window. You saw how to inspect the values of variables, insert breakpoints, and use the `Step`, `Go`, `Over` and `Out` buttons. You also got some practice debugging a function that didn’t work correctly using a four-step process for identifying and removing bugs:

1. Guess where the bug is located
2. Set a breakpoint and inspect the code
3. Identify the error and attempt to fix it
4. Repeat steps 1–3 until the error is fixed

Debugging is as much an art as it is a science. The only way to master debugging is to get a lot of practice with it!

One way to get some practice is to open the Debug Control Window and use it to step through your code as you work on the exercises and challenges throughout the rest of this book.

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-7

Additional Resources

For more information on debugging Python code, check out the following resources:

- [Python Debugging With Pdb](#)
- Recommended resources on [realpython.com](#)

[Leave feedback on this section »](#)

Chapter 8

Conditional Logic and Control Flow

Nearly all of the code you have seen in this book is **unconditional**. That is, the code does not make any choices. Every line of code is executed in the order that is written or that functions are called, with possible repetitions inside of loops.

In this chapter, you will learn how to write programs that perform different actions based on different conditions using **conditional logic**. Paired with functions and loops, conditional logic allows you to write complex programs that handle many different situations.

In this chapter, you will learn how to:

- Compare the values of two or more variables
- Write `if` statements to control the flow of your programs
- Handle errors with `try` and `except`
- Apply conditional logic to create simple simulations

Let's get started!

[Leave feedback on this section »](#)

8.1 Compare Values

Conditional logic is based on performing different actions depending on whether or not some expression, called a **conditional**, is true or false. This idea is not specific to computers. Humans use conditional logic all the time to make decisions.

For example, the legal age for purchasing alcoholic beverages in the United States is 21. The statement “If you are at least 21 years old, then you may purchase a beer” is an example of conditional logic. The phrase “you are at least 21 years old” is a conditional because it may be either true or false.

In computer programming, conditionals often take the form of comparing two values, such as determining if one value is greater than another, or whether or not two values are equal to each other. A standard set of symbols called **boolean comparators** are used to make comparisons, and most of them may already be familiar to you.

The following table describes these boolean comparators:

Boolean Comparator	Example	Meaning
>	a > b	a greater than b
<	a < b	a less than b
>=	a >= b	a greater than or equal to b
<=	a <= b	a less than or equal to b
!=	a != b	a not equal to b
==	a == b	a equal to b

The term **boolean** is derived from the last name of the English mathematician George Boole, whose works helped lay the foundations of modern computing. In Boole’s honor, conditional logic is sometimes called **boolean logic**, and conditionals are sometimes called **boolean expressions**.

There is also a fundamental data type called the **boolean**, or `bool` for short, which can have only one of two values. In Python, these values

are conveniently named `True` and `False`:

```
>>> type(True)
<class 'bool'>

>>> type(False)
<class 'bool'>
```

Note that `True` and `False` both start with capital letters.

The result of evaluating a conditional is always a boolean value:

```
>>> 1 == 1
True

>>> 3 > 5
False
```

In the first example, since `1` is equal to `1`, the result of `1 == 1` is `True`. In the second example, `3` is not greater than `5`, so the result is `False`.

Important

A common mistake when writing conditionals is to use the assignment operator `=`, instead of `==`, to test whether or not two values are equal.

Fortunately, Python will raise a `SyntaxError` if this mistake is encountered, so you'll know about it before you run your program.

You may find it helpful to think of boolean comparators as asking a question about two values. `a == b` asks whether or not `a` and `b` have the same value. Likewise, `a != b` asks whether or not `a` and `b` have different values.

Conditional expressions are not limited to comparing numbers. You may also compare values such as strings:

```
>>> "a" == "a"
```

```
True
```

```
>>> "a" == "b"
```

```
False
```

```
>>> "a" < "b"
```

```
True
```

```
>>> "a" > "b"
```

```
False
```

The last two examples above may look funny to you. How could one string be greater than or less than another?

The comparators `<` and `>` represent the notions of greater than and less than when used with numbers, but more generally they represent the notion of order. In this regard, `"a" < "b"` checks if the string `"a"` comes before the string `"b"`. But how are string ordered?

In Python, strings are ordered **lexicographically**, which is a fancy way to say they are ordered as they would appear in a dictionary. So you can think of `"a" < "b"` as asking whether or not the letter `a` comes before the letter `b` in the dictionary.

Lexicographic ordering extends to strings with two or more characters by looking at each component letter of the string:

```
>>> "apple" < "astronaut"
```

```
True
```

```
>>> "beauty" > "truth"
```

```
False
```

Since strings can contain characters other than letters of the alphabet, the ordering must extend to those other characters as well.

We won't go in to the details of how characters other than letters are

ordered. In practice, the `<` and `>` comparators are most often used with numbers, not strings.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. For each of the following conditional expressions, guess whether they evaluate to `True` or `False`. Then type them into the interactive window to check your answers:

```
1 <= 1
1 != 1
1 != 2
"good" != "bad"
"good" != "Good"
123 == "123"
```

2. For each of the following expressions, fill in the blank (indicated by `_`) with an appropriate boolean comparator so that the expression evaluates to `True`:

```
3 __ 4
10 __ 5
"jack" __ "jill"
42 __ "42"
```

[Leave feedback on this section »](#)

8.2 Add Some Logic

In addition to boolean comparators, Python has special keywords called **logical operators** that can be used to combine boolean expressions. There are three logical operators: `and`, `or`, and `not`.

Logical operators are used to construct compound logical expressions. For the most part, these have meanings similar to their meaning in the English language, although the rules regarding their use in Python are much more precise.

The and Keyword

Consider the following statements:

1. Cats have four legs.
2. Cats have tails.

In general, both of these statements are true.

When we combine these two statements using `and`, the resulting sentence “cats have four legs and cats have tails” is also a true statement. If both statements are negated, the compound statement “cats do not have four legs and cats do not have tails” is false.

Even when we mix and match false and true statements, the compound statement is false. “Cats have four legs and cats do not have tails” and “cats do not have four legs and cats have tails” are both false statements.

When two statements P and Q are combined with `and`, the **truth value** of the compound statement “ P and Q ” is true if and only if both P and Q are true.

Python’s `and` operator works exactly the same way. Here are four example of compound statements with `and`:

```
>>> 1 < 2 and 3 < 4 # Both are True  
True
```

Both statements are `True`, so the combination is also `True`.

```
>>> 2 < 1 and 4 < 3 # Both are False  
False
```

Both statements are `False`, so their combination is also `False`.

```
>>> 1 < 2 and 4 < 3 # Second statement is False  
False
```

`1 < 2` is `True`, but `4 < 3` is `False`, so their combination is `False`.

```
>>> 2 < 1 and 3 < 4 # First statement is False  
False
```

$2 < 1$ is False, and $3 < 4$ is True, so their combination is False.

The following table summarizes the rules for the `and` operator:

Combination using <code>and</code>	Result
True and True	True
True and False	False
False and True	False
False and False	False

You can test each of these rules in the interactive window:

```
>>> True and True  
True  
  
>>> True and False  
False  
  
>>> False and True  
False  
  
>>> False and False  
False
```

The `or` Keyword

When we use the word “or” in everyday conversation, sometimes we mean an **exclusive or**. That is, only the first option or the second option can be true.

For example, the phrase “I can stay or I can go” uses the exclusive or. I can’t both stay and go. Only one of these options can be true.

In Python the `or` keyword is inclusive. That is, if P and Q are two ex-

pressions, the statement “ P or Q ” is true if any of the following are true:

1. P is true
2. Q is true
3. Both P and Q are true

Let's look at some examples using numerical comparisons:

```
>>> 1 < 2 or 3 < 4  # Both are True
True

>>> 2 < 1 or 4 < 3  # Both are False
False

>>> 1 < 2 or 4 < 3  # Second statement is False
True

>>> 2 < 1 or 3 < 4  # First statement is False
True
```

Note that if any part of a compound statement is `True`, even if the other part is `False`, the result is always true `True`. The following table summarizes these results:

Combination using or	Result
True or True	True
True or False	True
False or True	True
False or False	False

Again, you can verify all of this in the interactive window:

```
>>> True or True
True
```

```
>>> True or False
```

```
True
```

```
>>> False or True
```

```
True
```

```
>>> False or False
```

```
False
```

The `not` Keyword

The `not` keyword reverses the truth value of a single expression:

Use of <code>not</code>	Result
<code>not True</code>	<code>False</code>
<code>not False</code>	<code>True</code>

You can verify this in the interactive window:

```
>>> not True
```

```
False
```

```
>>> not False
```

```
True
```

One thing to keep in mind with `not`, though, is that it doesn't always behave the way you might expect when combined with comparators like `==`. For example, `not True == False` returns `True`, but `False == not True` will raise an error:

```
>>> not True == False
```

```
True
```

```
>>> False == not True
```

```
File "<stdin>", line 1
```

```
False == not True
          ^
SyntaxError: invalid syntax
```

This happens because Python parses logical operators according to an **operator precedence**, just like arithmetic operators have an order of precedence in everyday math.

The order of precedence for logical and boolean operators, from highest to lowest, is described in the following table. Operators on the same row have equal precedence.

Operator Order of Precedence (Highest to Lowest)

<, <=, ==, >=, >
not
and
or

Looking again at the expression `False == not True`, `not` has a lower precedence than `==` in the order of operations. This means that when Python evaluates `False == not True`, it first tries to evaluate `False == not` which is syntactically incorrect.

You can avoid the `SyntaxError` by surrounding `not True` with parentheses:

```
>>> False == (not True)
True
```

Grouping expressions with parentheses is a great way to clarify which operators belong to which part of a compound expression.

Building Complex Expressions

You can combine the `and`, `or` and `not` keywords with `True` and `False` to create more complex expressions. Here's an example of a more complex expression:

```
True and not (1 != 1)
```

What do you think the value of this expression is?

To find out, break the expression down by starting on the far right side. `1 != 1` is `False`, since `1` has the same value as itself. So you can simplify the above expression as follows:

```
True and not (False)
```

`Now, not (False)` is the same as `not False`, which is `True`. So you can simplify the above expression once more:

```
True and True
```

Finally, `True and True` is just `True`. So, after a few steps, you can see that `True and not (1 != 1)` evaluates to `True`.

When working through complicated expressions, the best strategy is to start with the most complicated part of the expression and build outward from there.

For instance, try evaluating the following expression:

```
("A" != "A") or not (2 >= 3)
```

Start by evaluating the two expressions in parentheses. `"A" != "A"` is `False` because `"A"` is equal to itself. `2 >= 3` is also `False` because `2` is smaller than `3`. This gives you the following equivalent, but simpler, expression:

```
(False) or not (False)
```

Since `not` has a higher precedence than `or`, the above expression is equivalent to the following:

```
False or (not False)
```

`not False` is `True`, so you can simplify the expression once more:

```
False or True
```

Finally, since any compound expression with `or` is `True` if any one of the expressions on the left or right of the `or` is `True`, you can conclude that `("A" != "A") or not (2 >= 3)` is `True`.

Grouping expressions in a compound conditional statement with parentheses improves readability. Sometimes, though, parenthesis are required to produce the expected value.

For example, upon first inspection, you may expect the following to output `True`, but it actually returns `False`:

```
>>> True and False == True and False  
False
```

The reason this is `False` is that the `==` operator has a higher precedence than `and`, so Python interprets the expression as `True and (False == True) and False`. Since `False == True` is `False`, this is equivalent to `True and False and False`, which evaluates to `False`.

The following shows how to add parentheses so that the expression evaluates to `True`:

```
>>> (True and False) == (True and False)  
True
```

Logical operators and boolean comparators can be confusing the first time you encounter them, so if you don't feel like the material in this section comes naturally, don't worry!

With a little bit of practice, you'll be able to make sense of what's going on and build your own compound conditional statements when you need them.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

- Figure out what the result will be (True or False) when evaluating the following expressions, then type them into the interactive window to check your answers:

```
(1 <= 1) and (1 != 1)  
not (1 != 2)  
("good" != "bad") or False  
("good" != "Good") and not (1 == 1)
```

- Add parentheses where necessary so that each of the following expressions evaluates to True:

```
False == not True  
True and False == True and False  
not True and "A" == "B"
```

[Leave feedback on this section »](#)

8.3 Control the Flow of Your Program

Now that we can compare values to one other with boolean comparators and build complex conditional statements with logical operators, we can add some logic to our code so that it performs different actions for different conditions.

The if Statement

An if statement tells Python to only execute a portion of code if a condition is met.

For example, the following if statement will print 2 and 2 is 4 if the conditional `2 + 2 == 4` is True:

```
if 2 + 2 == 4:  
    print("2 and 2 is 4")
```

In English, you can read this as “if $2 + 2$ is 4, then print the string ‘2 and 2 is 4’.”

Just like while loops, an if statement has three parts:

1. The `if` keyword
2. A test condition, followed by a colon
3. An indented block of code that is executed if the test condition is `True`

In the above example, the test condition is `2 + 2 == 4`. Since this expression is `True`, executing the `if` statement in IDLE displays the text `2 and 2 is 4`.

If the test condition is `False` (for instance, `2 + 2 == 5`), Python skips over the indented block of code and continues execution on the next non-indented line.

For example, the following `if` statement does not print anything:

```
if 2 + 2 == 5:  
    print("Is this the mirror universe?")
```

A universe where `2 + 2 == 5` is `True` would be pretty strange indeed!

Note

Leaving off the colon (`:`) after the test condition in an `if` statement raises a `SyntaxError`:

```
>>> if 2 + 2 == 4  
SyntaxError: invalid syntax
```

Once the indented code block in an `if` statement is executed, Python will continue to execute the rest of the program.

Consider the following script:

```
grade = 95  
  
if grade >= 70:  
    print("You passed the class!")
```

```
print("Thank you for attending.")
```

The output looks like this:

```
You passed the class!  
Thank you for attending.
```

Since grade is 95, the test condition `grade >= 70` is True and the string "You passed the class!" is printed. Then the rest of the code is executed and "Thank you for attending." is printed.

If you change the value of grade to 40, the output looks like this:

```
Thank you for attending.
```

The line `print("Thank you for attending.")` is executed whether or not grade is greater than or equal to 70 because it is after the indented code block in the `if` statement.

A failing student will not know that they failed if all they see from your code is the text "Thank you for attending.".

Let's add another `if` statement to tell the student they did not pass if their grade is less than 70:

```
grade = 40  
  
if grade >= 70:  
    print("You passed the class!")  
  
if grade < 70:  
    print("You did not pass the class :(")  
  
print("Thank you for attending.")
```

The output now looks like this:

```
You did not pass the class :(  
Thank you for attending.
```

In English, we can describe an alternate case with the word “otherwise.” For instance, “If your grade is 70 or above, you pass the class. Otherwise, you do not pass the class.”

Fortunately, there is a keyword that does for Python what the word “otherwise” does in English.

The `else` Keyword

The `else` keyword is used after an `if` statement in order to execute some code only if the `if` statement’s test condition is `False`.

The following script uses `else` to shorten the code in the previous script for displaying whether or not a student passed a class:

```
grade = 40  
  
if grade >= 70:  
    print("You passed the class!")  
else:  
    print("You did not pass the class :(")  
  
print("Thank you for attending.")
```

In English, the `if` and `else` statements together read as ”If the grade is at least 70, then print the string ”You passed the class!”; otherwise, print the string ”You did not pass the class :(“.

Notice that the `else` keyword has no test condition, and is followed by a colon. No condition is needed, because it executes for any condition that fails the `if` statement’s test condition.

Important

Leaving off the colon (:) from the `else` keyword will raise a `SyntaxError`:

```
>>> if 2 + 2 == 5:  
...     print("Who broke my math?")  
... else  
SyntaxError: invalid syntax
```

The output from the above script is:

```
You did not pass the class :(  
Thank you for attending.
```

The line that prints "Thank you for attending." still runs, even if the indented block of code after `else` is executed.

The `if` and `else` keywords work together nicely if you only need to test a condition with exactly two states.

Sometimes, you need to check three or more conditions. For that, you use `elif`.

The `elif` Keyword

The `elif` keyword is short for “else if” and can be used to add additional conditions after an `if` statement.

Just like `if` statements, `elif` statements have three parts:

1. The `elif` keyword
2. A test condition, followed by a colon
3. An indented code block that is executed if the test condition evaluates to True

Important

Leaving off the colon (:) at the end of an `elif` statement raises a `SyntaxError`:

```
>>> if 2 + 2 == 5:  
...     print("Who broke my math?")  
... elif 2 + 2 == 4  
SyntaxError: invalid syntax
```

The following script combines `if`, `elif`, and `else` to print the letter grade a student earned in a class:

```
grade = 85 # 1  
  
if grade >= 90: # 2  
    print("You passed the class with a A.")  
elif grade >= 80: # 3  
    print("You passed the class with a B.")  
elif grade >= 70: # 4  
    print("You passed the class with a C.")  
else: # 5  
    print("You did not pass the class :(")  
  
print("Thanks for attending.") # 6
```

Both `grade >= 80` and `grade >= 70` are `True` when `grade` is 85, so you might expect both `elif` blocks on lines 3 and 4 to be executed.

However, only the first block for which the test condition is `True` is executed. All remaining `elif` and `else` blocks are skipped, so executing the script has the following output:

```
You passed the class with a B.  
Thanks for attending.
```

Let's break down the execution of the script step-by-step:

1. grade is assigned the value 85 in the line marked 1.
2. grade ≥ 90 is False, so the if statement marked 2 is skipped.
3. grade ≥ 80 is True, so the block under the elif statement in line 3 is executed, and "You passed the class with a B." is printed.
4. The elif and else statements in lines 4 and 5 are skipped, since the condition for the elif statement on line 3 was met.
5. Finally, line 6 is executed and "Thanks for attending." is printed.

The if, elif, and else keywords are some of the most commonly used keywords in the Python language. They allow you to write code that responds to different conditions with different behavior.

The if statement allows you to solve more complex problems than code without any conditional logic. You can even nest an if statement inside another one to write code that handles tremendously complex logic!

Nested if Statements

Just like for and while loops can be nested within one another, you nest an if statement inside another to create complicated decision making structures.

Consider the following scenario. Two people play a one-on-one sport against one another. You must decide which of two players wins depending on the players' scores and the sport they are playing:

- If the two players are playing basketball, the player with the greatest score wins.
- If the two players are playing golf, then the player with the lowest score wins.
- In either sport, if the two scores are equal, the game is a draw.

The following program solves this using nested if statements:

```
sport = input("Enter a sport: ")
p1_score = int(input("Enter player 1 score: "))
p2_score = int(input("Enter player 2 score: "))

# 1
if sport.lower() == "basketball":
    if p1_score == p2_score:
        print("The game is a draw.")
    elif p1_score > p2_score:
        print("Player 1 wins.")
    else:
        print("Player 2 wins.")

# 2
elif sport.lower() == "golf":
    if p1_score == p2_score:
        print("The game is a draw.")
    elif p1_score < p2_score:
        print("Player 1 wins.")
    else:
        print("Player 2 wins.")

# 3
else:
    print("Unknown sport")
```

This program first asks the user to input a sport and the scores for two players.

In (#1), the string assigned to `sport` is converted to lowercase using `.lower()` and is compared it to the string "basketball". This ensures that user input such as "Basketball" or "BasketBall" all get interpreted as the same sport.

Then the players scores are compared. If they are equal, the game is a draw. If player 1's score is larger than player 2's score, then player 1 wins the basketball game. Otherwise, player 2 wins the basketball game.

In (#2), the string assigned to `sport` is converted to lowercase gain and compared to the string "golf". Then the players scores are checked again. If the two scores are equal, the game is a draw. If player 1's score is less than player 2's score, then player 1 wins. Otherwise, player 2 wins.

Finally, in (#3), if the `sport` variable is assigned to a string other than "basketball" or "golf", the message "Unknown sport" is displayed.

The output of the script depends on the input value. Here's a sample execution using "basketball" as the sport:

```
Enter a sport: basketball
Player 1 score: 75
Player 2 score: 64
Player 1 wins.
```

Here's the output with the same player scores and the sport changed to "golf":

```
Enter a sport: golf
Player 1 score: 75
Player 2 score: 64
Player 2 wins.
```

If you enter anything besides basketball or golf for the sport, the program displays Unknown sport.

All together, there are seven possible ways that the program can run, which are described in the following table:

Sport	Score values
"basketball"	<code>p1_score == p2_score</code>
"basketball"	<code>p1_score > p2_score</code>
"basketball"	<code>p1_score < p2_score</code>
"golf"	<code>p1_score == p2_score</code>
"golf"	<code>p1_score > p2_score</code>
"golf"	<code>p1_score < p2_score</code>

Sport	Score values
everything else	any combination

Nested `if` statements can create many possible ways that your code can run. If you have many deeply nested `if` statements (more than two levels), then the number of possible ways the code can execute grows quickly.

Note

The complexity that results from using deeply nested `if` statements may make it difficult to predict how your program will behave under given conditions.

For this reason, nested `if` statements are generally discouraged.

Let's see how we simplify the previous program by removing nested `if` statements.

First, regardless of the sport, the game is a draw if `p1_score` is equal to `p2_score`. So, we can move the check for equality out from the nested `if` statements under each sport to make a single `if` statement:

```
if p1_score == p2_score:
    print("The game is a draw.")

elif sport.lower() == "basketball":
    if p1_score > p2_score:
        print("Player 1 wins.")
    else:
        print("Player 2 wins.")

elif sport.lower() == "golf":
    if p1_score < p2_score:
        print("Player 1 wins.")
    else:
```

```
    print("Player 2 wins.")
```

```
else:  
    print("Unknown sport.")
```

Now there are only six ways that the program can execute.

That's still quite a few ways. Can you think of any way to make the program simpler?

Here's one way to simplify it. Player 1 wins if the sport is basketball and their score is greater than player 2's score, or if the sport is golf and their score is less than player 2's score.

We can describe this with compound conditional expressions:

```
sport = sport.lower()  
p1_wins_basketball = (sport == "basketball") and (p1_score > p2_score)  
p1_wins_golf = (sport == "golf") and (p1_score < p2_score)  
p1_wins = player1_wins_basketball or player1_wins_golf
```

This code is pretty dense, so let's walk through it one step at a time.

First the string assigned to `sport` is converted to all lowercase so that we can compare the value to other strings without worrying about errors due to case.

On the next line, we have a structure that might look a little strange. There is an assignment operator (`=`) followed by an expression with the equality comparator (`==`). This line evaluates the following compound logical expression and assigns its value to the `p1_wins_basketball` variable:

```
(sport == "basketball") and (p1_score > p2_score)
```

If `sport` is "basketball" and player 1's score is larger than player 2's score, then `p1_wins_basketball` is True.

Next, a similar operation is done for the `p1_wins_golf` variable. If score

is "golf" and player 1's score is less than player 2's score, then `p1_wins_golf` is `True`.

Finally, `p1_wins` will be `True` if player 1 wins the basketball game or the golf game, and will be `False` otherwise.

Using this code, you can simplify the program quite a bit:

```
if p1_score == p2_score:
    print("The game is a draw.")

elif (sport.lower() == "basketball") or (sport.lower() == "golf"):
    sport = sport.lower()
    p1_wins_basketball = (sport == "basketball") and (p1_score > p2_score)
    p1_wins_golf = (sport == "golf") and (p1_score < p2_score)
    p1_wins = p1_wins_basketball or p1_wins_golf

    if p1_wins:
        print("Player 1 wins.")
    else:
        print("Player 2 wins.")

else:
    print("Unknown sport")
```

In this revised version of the program, there are only four ways the program can execute, and the code is easier to understand.

Nested `if` statements are sometimes necessary. However, if you find yourself writing lots of nested `if` statements, it might be a good idea to stop and think about how you might simplify your code.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Write a script that prompts the user to enter a word using the `input()` function, stores that input in a variable, and then displays whether the length of that string is less than 5 characters, greater than 5 characters, or equal to 5 characters by using a set of `if`, `elif`

and `else` statements.

[Leave feedback on this section »](#)

8.4 Challenge: Find the Factors of a Number

A factor of a positive integer n is any positive integer less than or equal to n that divides n with no remainder.

For example, 3 is a factor of 12 because 12 divided by 3 is 4, with no remainder. However, 5 is not a factor of 12 because 5 goes into 12 twice with a remainder of 2.

Write a script `factors.py` that asks the user to input a positive integer and then prints out the factors of that number. Here's a sample run of the program with output:

```
Enter a positive integer: 12
1 is a factor of 12
2 is a factor of 12
3 is a factor of 12
4 is a factor of 12
6 is a factor of 12
12 is a factor of 12
```

Hint: Recall from Chapter 5 that you can use the `%` operator to get the remainder of dividing one number by another.

You can find the solutions to this code challenge and many other bonus resources online at realpython.com/python-basics/resources.

[Leave feedback on this section »](#)

8.5 Break Out of the Pattern

In Chapter 6 you learned how to repeat a block of code many times using a `for` or `while` loop. Loops are useful for performing a repetitive task and for applying some processing to many different inputs.

Combining `if` statements with `for` loops opens up powerful techniques for controlling how code is run.

In this section, you'll learn how to write `if` statements that are nested in `for` loops and learn about two keywords — `break` and `continue` — that allow you to more precisely control the flow of execution through a loop.

if Statements and for Loops

The block of code in a `for` loop is just like any other block of code. That means you can nest an `if` statement in a `for` loop just like you can anywhere else in your code.

The following example uses a `for` loop with an `if` statement to compute and display the sum of all even integers less than 100:

```
sum_of_evens = 0

for n in range(1, 100):
    if n % 2 == 0:
        sum_of_evens = sum_of_evens + n

print(sum_of_evens)
```

First, the `sum_of_evens` variable is initialized to 0. Then the program loops over the numbers 1 to 99 and adds the even values to `sum_of_evens`. The final value of `sum_of_evens` is 2450.

break

The `break` keyword tells Python to literally break out of a loop. That is, the loop stops completely and any code after the loop is executed.

For example, the following code loops over the numbers 0 to 3, but stops the loop when the number 2 is encountered:

```
for n in range(0, 4):
    if n == 2:
        break
    print(n)

print(f"Finished with n = {n}")
```

Only the first two numbers are printed in the output:

```
0
1
Finished with n = 2
```

continue

The `continue` keyword is used to skip any remaining code in the loop body and continue on to the next iteration.

For example, the following code loops over the numbers 0 to 3, printing each number as is goes, but skips the number 2:

```
for i in range(0, 4):
    if i == 2:
        continue
    print(i)

print(f"Finished with i = {i}")
```

All the numbers except for 2 are printed in the output:

```
0  
1  
3  
Finished with i = 3
```

Note

It's always a good idea to give short but descriptive names to your variables that make it easy to tell what they are supposed to represent.

The letters `i`, `j` and `k` are exceptions because they are so common in programming.

These letters are almost always used when we need a “throw-away” number solely for the purpose of keeping count while working through a loop.

To summarize, the `break` keyword is used to stop a loop if a certain condition is met, and the `continue` keyword is used to skip an iteration of a loop when a certain condition is met.

for...else Loops

Loops can have their own `else` clause in Python, although this structure isn't used very frequently.

Let's look at an example:

```
phrase = "it marks the spot"

for character in phrase:
    if character == "X":
        break
    else:
        print("There was no 'X' in the phrase")
```

The `for` loop in this example loops over the characters phrase "it marks

the spot" and stops if the letter "X" is found.

If you run the code in the example, you'll see that There was no 'X' in the phrase is printed to the console.

Now try changing phrase to the string "X marks the spot". When you run the same code with this phrase, there is no output. What's going on?

Any code in the else block after a for loop is executed only if the for loop completes without encountering a break statement.

So, when you run the code with phrase = "it marks the spot", the line of code containing the break statement is never run since there is no x character in the phrase, which means that the else block is executed and the string "There was no 'X' in the phrase" is displayed.

On the other hand, when you run the code with phrase = "X marks the spot", the line containing the break statement *does* get executed, so the else block is never run and no output gets displayed.

Here's a practical example that gives a user three attempts to enter a password:

```
for n in range(3):
    password = input("Password: ")
    if password == "I<3Bieber":
        break
    print("Password is incorrect.")
else:
    print("Suspicious activity. The authorities have been alerted.")
```

This example loops over the number 0 to 2. On each iteration, the user is prompted to enter a password. If the password entered is correct, then break is used to exit the loop. Otherwise, the user is told that the password is incorrect and given another attempt.

After three unsuccessful attempts, the for loop terminates without ever executing the line of code containing break. In that case, the else

block is executed and the user is warned that the authorities have been alerted.

Note

We have focused on `for` loops in this section because they are generally the most common kind of loops.

However, everything discussed here also works for `while` loops. That is, you can use `break` and `continue` inside a `while` loop. `while` loops can even have an `else` clause!

Using conditional logic inside the body of a loop opens up several possibilities for controlling how your code executes.

You can stop loops early with the `break` keyword or skip an iteration with `continue`. You can even make sure some code only runs if a loop completes without ever encountering a `break` statement.

These are some powerful tools to have in your tool kit!

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Using `break`, write a program that repeatedly asks the user for some input and only quits if the user enters "q" or "Q".
2. Using `continue`, write a program that loops over the numbers 1 to 50 and prints all numbers that are not multiples of 3.

[Leave feedback on this section »](#)

8.6 Recover From Errors

Encountering errors in your code might be frustrating, but it's totally normal! It happens to even the best programmers.

Programmers often refer to run-time errors as **exceptions**. So, when you encounter an error, congratulate yourself! You've just made the code do something exceptional!

Errors aren't always a bad thing. That is, they don't always mean you made a mistake. For example, trying to divide the 1 by 0 results in a `ZeroDivisionError`. If the divisor is entered by a user, you have no way of knowing ahead of time whether or not the user will enter a 0!

In order to create robust programs, you need to be able to handle errors caused by invalid user input — or any other unpredictable source. In this section you'll learn how.

A Zoo of Exceptions

When you encounter an exception, it's useful to know what went wrong. Python has a number of built-in exception types that describe different kinds of errors.

Throughout this book you have seen several different errors. Let's collect them here and add a few new ones to the list.

ValueError

A `ValueError` occurs when an operation encounters an invalid value. For example, trying to convert the string "not a number" to an integer results in a `ValueError`:

```
>>> int("not a number")
Traceback (most recent call last):
  File "<pyshell#1>", line 1, in <module>
    int("not a number")
ValueError: invalid literal for int() with base 10: 'not a number'
```

The name of the exception is displayed on the last line, followed by a description of the specific problem that occurred. This is the general format for all Python exceptions.

TypeError

A `TypeError` occurs when an operation is performed on a value of the wrong type. For example, trying to add a string and an integer will result in a `TypeError`:

```
>>> "1" + 2
Traceback (most recent call last):
  File "<pyshell#1>", line 1, in <module>
    "1" + 2
TypeError: can only concatenate str (not "int") to str
```

NameError

A `NameError` occurs when you try to use a variable name that hasn't been defined yet:

```
>>> print(does_not_exist)
Traceback (most recent call last):
  File "<pyshell#3>", line 1, in <module>
    print(does_not_exist)
NameError: name 'does_not_exist' is not defined
```

ZeroDivisionError

A `ZeroDivisionError` occurs when the divisor in a division operation is 0:

```
>>> 1 / 0
Traceback (most recent call last):
  File "<pyshell#4>", line 1, in <module>
    1 / 0
ZeroDivisionError: division by zero
```

OverflowError

An `OverflowError` occurs when the result of an arithmetic operation is too large. For example, trying to raise the value 2.0 to the power 1_000_000 results in an `OverflowError`:

```
>>> pow(2.0, 1_000_000)
Traceback (most recent call last):
  File "<pyshell#6>", line 1, in <module>
    pow(2.0, 1_000_000)
OverflowError: (34, 'Result too large')
```

You may recall from Chapter 5 that integers in Python have unlimited precision. This means that `OverflowErrors` can only occur with floating-point numbers.

Raising the `integer` 2 to the value `1_000_000` will not raise an `OverflowError`!

A list of Python's built-in exceptions can be found [in the docs](#).

The `try` and `except` Keywords

Sometimes you can predict that a certain exception might occur. Instead of letting the program crash, you can catch the error if it occurs and do something else instead.

For example, you might need to ask the user to input an integer. If the user enters a non-integer value, such as the string "a", you need to let them know that they entered an invalid value.

To prevent the program from crashing you can use the `try` and `except` keywords. Let's look at an example:

```
try:
    number = int(input("Enter an integer: "))
except ValueError:
    print("That was not an integer")
```

The `try` keyword is used to indicate a `try` block and is followed by a colon. The code indented after `try` is executed. In this case, the user is asked to input an integer. Since `input()` returns a string, the user input is converted to an integer with `int()` and the result is assigned to the variable `number`.

If the user inputs a non-integer value, the `int()` operation will raise a `ValueError`. If that happens, the code indented below the line `except ValueError` is executed. So, instead of the program crashing, the string "That was not an integer" is displayed.

If the user does input a valid integer value, then the code in the `except ValueError` block is never executed.

On the other hand, if a different kind of exception had occurred, such as a `TypeError`, then the program will crash. The above example only handles one type of exception — a `ValueError`.

You can handle multiple exception types by separating the exception names with commas and putting the list of names in parentheses:

```
def divide(num1, num2):
    try:
        print(num1 / num2)
    except (TypeError, ZeroDivisionError):
        print("encountered a problem")
```

In this example, the function `divide()` takes two parameters `num1` and `num2` and prints the result of dividing `num1` by `num2`.

If `divide()` is called with an argument that is a string, then the division operation will raise a `TypeError`. Additionally, if `num2` is 0, then a `ZeroDivisionError` is raised.

The line `except (TypeError, ZeroDivisionError)` will handle both of these exceptions and display the string "encountered a problem" if either exception is raised.

Many times, though, it is helpful to catch each error individually so that you can display text that is more helpful to the user. To do this, you can use multiple `except` blocks after a `try` block:

```
def divide(num1, num2):
    try:
```

```
print(num1 / num2)
except TypeError:
    print("Both arguments must be numbers")
except ZeroDivisionError:
    print("num2 must not be 0")
```

In this example, the `ValueError` and `ZeroDivisionError` are handled separately. This way, a more descriptive message is displayed if something goes wrong.

If one of `num1` or `num2` is not a number, then a `TypeError` is raised and the message "Both arguments must be numbers" is displayed. If `num2` is 0, then a `ZeroDivisionError` is raised and the message "num2 must not be 0" is displayed.

The “Bare” `except` Clause

You can use the `except` keyword by itself without naming specific exceptions:

```
try:
    # Do lots of hazardous things that might break
except:
    print("Something bad happened!")
```

If any exception is raised while executing the code in the `try` block, the `except` block will run and the message "Something bad happened!" will be displayed.

This might sound like a great way to ensure your program never crashes, **but this is actually bad idea and the pattern is generally frowned upon!**

There are a couple of reasons for this, but the most important reason for new programmers is that catching every exception could hide bugs in your code, giving you a false sense of confidence that your code works as expected.

If you only catch specific exceptions, then when unexpected errors are encountered, Python will print the traceback and error information giving you more information to work with when debugging your code.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Write a script that repeatedly asks the user to input an integer, displaying a message to “try again” by catching the `ValueError` that is raised if the user did not enter an integer.

Once the user enters an integer, the program should display the number back to the user and end without crashing.

2. Write a program that asks the user to input a string and an integer n . Then display the character at index n in the string.

Use error handling to make sure the program doesn’t crash if the user does not enter an integer or the index is out of bounds. The program should display a different message depending on what error occurs.

[Leave feedback on this section »](#)

8.7 Simulate Events and Calculate Probabilities

In this section, we’ll apply some of the concepts we’ve learned about loops and conditional logic to a real world problem: simulating events and calculating probabilities.

We’ll be running a simple simulation known as a **Monte Carlo** experiment. Each experiment consists of a **trial**, which is just some process that can be repeated — such as flipping a coin — that generated some outcome — such as landing on heads or tails. The trial is repeated over

and over again in order to calculate the probability that some outcome occurs.

In order to do this, we need to add some randomness to our code.

The `random` module

Python provides several functions for generating random numbers in the `random` module. A **module** is a collection of related code. Python's **standard library** is an organized collection of modules that you can **import** into your own code in order to solve various problems.

To import the `random` module, type the following into IDLE's interactive window:

```
>>> import random
```

Now we can use functions from the `random` module in our code. For example, the `randint()` has two required parameters called `a` and `b` and returns a random integer that is greater than or equal to `a` and less than or equal to `b`. Both `a` and `b` must be integers.

For example, the following code produces a random integer between 1 and 10:

```
>>> random.randint(1, 10)
9
```

Since the result is random, your output will probably be different than 9. If you type the same code in again, you will likely get a different number.

Since `randint()` is located in the `random` module, you must type `random` followed by a dot (.) and then the function name in order to use it.

It is important to remember that when using `randint()`, the two parameters `a` and `b` must both be integers, and the output might be equal to one of `a` and `b`, or any number in-between. For instance, `random.randint(0, 1)` randomly returns either a 0 or a 1.

Furthermore, each integer between `a` and `b` is equally likely to be returned by `randint()`. So, for `randint(1, 10)`, each integer between 1 and 10 has a 10% chance of being returned. For `randint(0, 1)`, there is a 50% chance a 0 is returned.

Flipping Fair Coins

Let's see how to use `randint()` to simulate flipping a fair coin. By a fair coin, we mean a coin that, when flipped, has an equal chance of landing on heads or tails.

One trial for our experiment will be flipping the coin. The outcome is either a heads or a tails. The question is: in general, over many coin flips, what is the ratio of heads to tails?

Let's think about how to solve this problem. We'll need to keep track of how many times we get a heads or tails, so we need a heads tally and a tails tally. Each trial has two steps:

1. Flip the coin.
2. If the coin lands on heads, update the heads tally. Otherwise, the coin lands on tails so update the tails tally.

We need to repeat the trial many times, say 10,000. A `for` loop over `range(10_000)` is a good choice for doing something like that.

Now that we have a plan, let's start by writing a function called `coin_flip()` that randomly returns the string "heads" or the string "tails". We can do this using `random.randint(0, 1)`. We'll use 0 to represent heads and 1 for tails.

Here's the code for the `coin_flip()` function:

```
import random

def coin_flip():
    """Randomly return 'heads' or 'tails'."""


```

```
if random.randint(0, 1) == 0:  
    return "heads"  
else:  
    return "tails"
```

If `random.randint(0, 1)` returns a 0, then `coin_flip()` returns "heads". Otherwise, `coin_flip()` returns "tails".

Now we can write a `for` loop that flips the coin 10,000 times and updates a heads or tails tally accordingly:

```
# First initialize the tallies to 0  
heads_tally = 0  
tails_tally = 0  
  
for trial in range(10_000):  
    if coin_flip() == "heads":  
        heads_tally = heads_tally + 1  
    else:  
        tails_tally = tails_tally + 1
```

First, two variables `heads_tally` and `tails_tally` are created and both are initialized to the integer 0.

Then the `for` loop runs 10,000 times. Each time, the `coin_flip()` function is called. If it returns the string "heads", then the `heads_tally` variable is incremented by 1. Otherwise `tails_tally` is incremented by 1.

Finally, we can print the ratio of heads and tails:

```
ratio = heads_tally / tails_tally  
print(f"The ratio of heads to tails is {ratio}")
```

If you save the above code to a script and run it a few times, you will see that the result is usually between .98 and 1.02. If you increase the `range(10_000)` in the `for` loop to, say, `range(50_000)`, the results should get closer to 1.0.

This behavior makes sense. Since the coin is fair, we should expect

that after many flips, the number of heads is roughly equal to the number of tails.

In life, things aren't always fair. A coin may have a slight tendency to land on heads instead of tails, or vice versa. So, how do you simulate something like an unfair coin?

Tossing Unfair Coins

`randint()` returns a 0 or a 1 with equal probability. If 0 represents tails and 1 represents heads, then to simulate an unfair coin we need a way to return one of 0 or 1 with a higher probability.

The `random()` function can be called without any arguments and returns a floating-point number greater than or equal to 0.0 but less than 1.0. Each possible return value is equally likely. In probability theory, this is known as a [uniform probability distribution](#).

One consequence of this is that, given a number `n` between 0 and 1, the probability that `random()` returns a number less than `n` is just `n` itself. For example, the probability that `random()` is less than .8 is .8 and the probability that `random()` is less than .25 is .25.

Using this fact, we can write a function that simulates a coin flip, but returns tails with a specified probability:

```
import random

def unfair_coin_flip(probability_of_tails):
    if random.random() < probability_of_tails:
        return "tails"
    else:
        return "heads"
```

For example, `unfair_coin_flip(.7)` has a 70% chance of returning "tails".

Let's re-write the coin flip experiment from earlier using `unfair_coin_-`

`flip()` to run each trial with an unfair coin:

```
heads_tally = 0
tails_tally= 0

for trial in range(10_000):
    if unfair_coin_flip(.7) == "heads":
        heads_tally = heads_tally + 1
    else:
        tails_tally = tails_tally + 1

ratio = heads_tally / tails_tally
print(f"The ratio of heads to tails is {ratio}")
```

Running this simulation a few times shows that the ratio of heads to tails has gone down from 1 in the experiment with a fair coin to about .43.

In this section you learned about the `randint()` and `random()` functions in the `random` module and saw how to use conditional logic and loops to write some coin toss simulations. Simulations like these are used in numerous disciplines to make predictions and test computer models of real world events.

The `random` module provides many useful functions for generating random numbers and writing simulations. You can learn more about `random` in Real Python's [Generating Random Data in Python \(Guide\)](#).

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Write a function called `roll()` that uses the `randint()` function to simulate rolling a fair die by returning a random integer between 1 and 6.
2. Write a script that simulates 10,000 rolls of a fair die and displays the average number rolled.

[Leave feedback on this section »](#)

8.8 Challenge: Simulate a Coin Toss Experiment

Suppose you flip a fair coin repeatedly until it lands on both heads and tails at least once each. In other words, after the first flip, you continue to flip the coin until it lands on something different.

Doing this generates a sequence of heads and tails. For example, the first time you do this experiment, the sequence might be heads, heads, then tails.

On average, how many flips are needed for the sequence to contain both heads and tails?

Write a simulation that runs 10,000 trials of the experiment and prints the average number of flips per trial.

You can find the solutions to this code challenge and many other bonus resources online at realpython.com/python-basics/resources.

[Leave feedback on this section »](#)

8.9 Challenge: Simulate an Election

With some help from the `random` module and a little condition logic, you can simulate an election between two candidates.

Suppose two candidates, Candidate A and Candidate B, are running for mayor in a city with three voting regions. The most recent polls show that Candidate A has the following chances for winning in each region:

- Region 1: 87% chance of winning
- Region 2: 65% chance of winning

- Region 3: 17% chance of winning

Write a program that simulates the election 10,000 times and prints the percentage of where Candidate A wins.

To keep things simple, assume that a candidate wins the election if they win in at least two of the three regions.

You can find the solutions to this code challenge and many other bonus resources online at realpython.com/python-basics/resources.

[Leave feedback on this section »](#)

8.10 Summary and Additional Resources

In this chapter, you learned about conditional statements and conditional logic. You saw how to compare values using comparison operators like `<`, `>`, `<=`, `>=`, `!=`, and `==`. You also saw how to build complex conditional statements using `and`, `or` and `not`.

Next, you saw how to control the flow of your program using `if` statements. You learned how to create branches in your program using `if...else` and `if...elif...else`. You also learned how to control precisely how code is executed inside of an `if` block using `break` and `continue`.

You learned about the `try...except` pattern to handle errors that may occur during run-time. This is an important construct that allows your programs to handle the unexpected gracefully, and keep users of your programs happy that the program doesn't crash.

Finally, you applied the techniques you learned in this chapter and used the `random` module to build some simple simulations.

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-8

Additional Resources

A wise Vulcan once said:

Logic is the beginning of wisdom, not the end.

— Spock, *Star Trek*

Check out the following resources to learn more about conditional logic:

- [Operators and Expressions in Python](#)
- [Conditional Statements in Python](#)
- Recommended resources on realpython.com

[Leave feedback on this section »](#)

Chapter 9

Tuples, Lists, and Dictionaries

So far, you have been working with fundamental data types like `str`, `int`, and `float`. Many real-world problems are easier to solve when simple data types are combined into more complex data structures.

A **data structure** models a collection of data, such as a list of numbers, a row in a spreadsheet, or a record in a database. Modeling the data that your program interacts with using the right data structure is often the key to writing simple and effective code.

Python has three built-in data structures that are the focus of this chapter: tuples, lists, and dictionaries.

In this chapter, you will learn:

- How to work with tuples, lists, and dictionaries
- What immutability is and why it is important
- When to use different data structures

Let's dive in!

[Leave feedback on this section »](#)

9.1 Tuples Are Immutable Sequences

Perhaps the simplest compound data structure is a sequence of items.

A **sequence** is an ordered list of values. Each element in a sequence is assigned an integer, called an **index**, that determines the order in which the values appear. Just like strings, the index of the first value in a sequence is 0.

For example, the letters of the English alphabet form a sequence whose first element is A and last element is Z. Strings are also sequences. The string "Python" has six elements, starting with "P" at index 0, and "n" at index 5.

Some real-world examples of sequences include the values emitted by a sensor every second, the sequence of a student's test scores, or the sequence of daily stock values for some company over a period of time.

In this section, you'll learn how to use Python's built-in `tuple` data type to create sequences of values.

What is a Tuple?

The word **tuple** comes from mathematics, where it is used to describe a finite ordered sequence of values.

Usually, mathematicians write tuples by listing each element, separated by a comma, inside a pair of parentheses. `(1, 2, 3)` is a tuple containing three integers.

Tuples are **ordered** because their elements appear in an ordered fashion. The first element of `(1, 2, 3)` is 1, the second element is 2, and the third is 3.

Python borrows both the name and the notation for tuples from mathematics.

How to Create a Tuple

There are a few ways to create a tuple in Python. We will cover two of them:

1. Tuple literals
2. The `tuple()` built-in

Tuple Literals

Just like a string literal is a string that is explicitly created by surrounding some text with quotes, a **tuple literal** is a tuple that is written out explicitly as a comma-separated list of values surrounded by parentheses.

Here's an example of a tuple literal:

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

This creates a tuple containing the integers 1, 2, and 3, and assigns it to the name `my_first_tuple`.

You can check that `my_first_tuple` is a tuple using `type()`:

```
>>> type(my_first_tuple)
<class 'tuple'>
```

Unlike strings, which are sequences of characters, tuples may contain any type of value, including values of different types. The tuple `(1, 2.0, "three")` is perfectly valid.

There is a special tuple that doesn't contain any values. This tuple is called the **empty tuple** and can be created by typing two parentheses without anything between them:

```
>>> empty_tuple = ()
```

At first glance, the empty tuple may seem like a strange and useless concept, but it is actually quite practical.

For example, suppose you are asked to provide a tuple containing all the integers that are both even and odd. No such integer exists, but the empty tuple allows you to provide the requested tuple.

How do you think you create a tuple with exactly one element? Try out the following in IDLE:

```
>>> x = (1)
>>> type(x)
<class 'int'>
```

When you surround a value with parentheses, but don't include any commas, Python interprets the value not as a tuple but as the type of value inside the parentheses. So, in this case, (1) is just a weird way of writing the integer 1.

To create the tuple containing the single value 1, you need to include a comma after the 1:

```
>>> x = (1,)
>>> type(x)
<class 'tuple'>
```

A tuple containing a single element might seem as strange as the empty tuple. Couldn't you just drop all this tuple business and just use the value itself?

It all depends on the problem you are solving.

If you are asked to provide a tuple containing all prime numbers that are also even, you must provide the tuple (2,) since 2 is the only even prime number. The value 2 isn't a good solution because it isn't a tuple.

This might seem overly pedantic, but programming often involves a certain amount of pedantry. Computers are, after all, the ultimate pedants.

The `tuple()` Built-In

You can also use the `tuple()` built-in to create a tuple from another sequence type, such as a string:

```
>>> tuple("Python")
('P', 'y', 't', 'h', 'o', 'n')
```

`tuple()` only accepts a single parameter, so you can't just list the values you want in the tuple as individual arguments. If you do, Python raises a `TypeError`:

```
>>> tuple(1, 2, 3)
Traceback (most recent call last):
  File "<pyshell#0>", line 1, in <module>
    tuple(1, 2, 3)
TypeError: tuple expected at most 1 arguments, got 3
```

You will also get a `TypeError` if the argument passed to `tuple()` can't be interpreted as a list of values:

```
>>> tuple(1)
Traceback (most recent call last):
  File "<pyshell#1>", line 1, in <module>
    tuple(1)
TypeError: 'int' object is not iterable
```

The word **iterable** in the error message indicates that a single integer can't be **iterated**, which is to say that the integer data type doesn't contain multiple values that can be accessed one-by-one.

The single parameter of `tuple()` is optional, though, and leaving it out produces an empty tuple:

```
>>> tuple()
()
```

However, most Python programmers prefer to use the shorter `()` for creating an empty tuple.

Similarities Between Tuples and Strings

Tuples and strings have a lot in common. Both are sequence types with a finite length, support indexing and slicing, are immutable, and can be iterated over in a loop.

The main difference between strings and tuples is that the elements of tuples can be any kind of value you like, whereas strings can only contain characters.

Let's look at some of the parallels between strings in tuples in more depth.

Tuples Have a Length

Both strings and tuples have a **length**. The length of a string is the number of characters in it. The length of a tuple is the number of elements it contains.

Just like strings, the `len()` function can be used to determine the length of a tuple:

```
>>> numbers = (1, 2, 3)
>>> len(numbers)
3
```

Tuples Support Indexing and Slicing

Recall from Chapter 4 that you can access a character in a string using index notation:

```
>>> name = "David"
>>> name[1]
'a'
```

The index notation `[1]` after the variable `name` tells Python to get the character at index 1 in the string "David". Since counting starts at 0, the character at index 1 is the letter "a".

Tuples also support index notation:

```
>>> values = (1, 3, 5, 7, 9)
>>> values[2]
5
```

Another feature that strings and tuples have in common is slicing. Recall that you can extract a substring from a string using slicing notation:

```
>>> name = "David"
>>> name[2:4]
"vi"
```

The slice notation [2:4] after the variable `name` creates a new string containing the characters in `name` starting at position 2 and up to, but not including, the character at position 4.

Slicing notation also works with tuples:

```
>>> values = (1, 3, 5, 7, 9)
>>> values[2:4]
(5, 7)
```

The slice `values[2:4]` creates a new tuple containing the all integers in `values` starting at position 2 and up to, but not including, the integer at position 4.

The same rules governing string slices also apply to tuple slices. You may want to take some time to review the slicing examples in Chapter 4 with some of your own examples of tuples.

Tuples Are Immutable

Like strings, tuples are immutable. This means you can't change the value of an element of a tuple once it has been created.

If you do try to change the value at some index of a tuple, Python will raise a `TypeError`:

```
>>> values[0] = 2
Traceback (most recent call last):
  File "<pyshell#1>", line 1, in <module>
    values[0] = 2
TypeError: 'tuple' object does not support item assignment
```

Note

Although tuples are immutable, there are some situations in which the values in a tuple can change.

These quirks and oddities are covered in depth in Real Python's [Immutability in Python](#) video course.

Tuples Are Iterable

Just like strings, tuples are **iterable**, so you can loop over them:

```
>>> vowels = ("a", "e", "i", "o", "u")
>>> for vowel in vowels:
...     print(vowel.upper())
...
A
E
I
O
U
```

The `for` loop in this example works just like the `for` loops you saw in Chapter 6 that loop over a `range()` of numbers.

On the first step of the loop, the value "a" is extracted from the tuple `vowels`. It is converted to an upper case letter using the `.upper()` string method you learned about in Chapter 4, and then displayed with `print()`.

The next step of the loop extracts the value "e", converts it to upper case, and prints it. This continues for each of the values "i", "o", and

"u".

Now that you've seen how to create tuples and some of the basic operations they support, let's look at some common use cases.

Tuple Packing and Unpacking

There is a third, although less common, way of creating a tuple. You can type out a comma-separated list of values and leave off the parentheses:

```
>>> coordinates = 4.21, 9.29
>>> type(coordinates)
<class 'tuple'>
```

It looks like two values are being assigned to the single variable `coordinates`. In a sense, they are, although the result is that both values are **packed** into a single tuple. You can verify that `coordinates` is indeed a tuple with `type()`.

If you can pack values into a tuple, it only makes sense that you can unpack them as well:

```
>>> x, y = coordinates
>>> x
4.21
>>> y
9.29
```

Here the values contained in the single tuple `coordinates` are **unpacked** into two distinct variables `x` and `y`.

By combining tuple packing and unpacking, you can make multiple variable assignments in a single line:

```
>>> name, age, occupation = "David", 34, "programmer"
>>> name
'David'
```

```
>>> age  
34  
>>> occupation  
'programmer'
```

This works because first, on the right hand side of the assignment, the values "David", 34, and "programmer" are packed into a tuple. Then the values are unpacked into the three variables name, age, and programmer, in that order.

Note

While assigning multiple variables in a single line can shorten the number of lines in a program, you may want to refrain from assigning too many values in a single line.

Assigning more than two or three variables this way can make it difficult to tell which value is assigned to which variable name.

Keep in mind that the number of variable names on the left of the assignment expression must equal the number of values in the tuple on the right hand side, otherwise Python will raise a `ValueError`:

```
>>> a, b, c, d = 1, 2, 3  
Traceback (most recent call last):  
  File "<pyshell#0>", line 1, in <module>  
    a, b, c, d = 1, 2, 3  
ValueError: not enough values to unpack (expected 4, got 3)
```

The error message here tells you that the tuple on the right hand side doesn't have enough values to unpack into the four variable names.

Python also raises a `ValueError` if the number of values in the tuple exceeds the number of variable names:

```
>>> a, b, c = 1, 2, 3, 4  
Traceback (most recent call last):  
  File "<pyshell#1>", line 1, in <module>
```

```
a, b, c = 1, 2, 3, 4
```

```
ValueError: too many values to unpack (expected 3)
```

Now the error message indicates that there are too many values in the tuple to unpack into three variables.

Checking Existence of Values With `in`

You can check whether or not a value is contained in a tuple with the `in` keyword.

```
>>> vowels = ("a", "e", "i", "o", "u")
>>> "o" in vowels
True
>>> "x" in vowels
False
```

If the value to the left of `in` is contained in the tuple to the right of `in`, the result is `True`. Otherwise, the result is `False`.

Returning Multiple Values From a Function

One common use of tuples is to return multiple values from a single function.

```
>>> def adder_subtractor(num1, num2):
...     return (num1 + num2, num1 - num2)
...
>>> adder_subtractor(3, 2)
(5, 1)
```

The function `adder_subtractor()` has two parameters, `num1` and `num2`, and returns a tuple whose first element is the sum of the two numbers, and whose second element is the difference.

Strings and tuples are just two of Python's built-in sequence types. Both are immutable and iterable and can be used with index and slicing notation.

In the next section, you'll learn about a third sequence type with one very big difference from strings and tuples: mutability.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Create a tuple literal named `cardinal_numbers` that holds the strings "first", "second" and "third", in that order.
2. Using index notation and `print()`, display the string at index 1 in `cardinal_numbers`.
3. Unpack the values in `cardinal_numbers` into three new strings named `position1`, `position2` and `position3` in a single line of code, then print each value on a separate line.
4. Create a tuple called `my_name` that contains the letters of your name by using `tuple()` and a string literal.
5. Check whether or not the character "x" is in `my_name` using the `in` keyword.
6. Create a new tuple containing all but the first letter in `my_name` using slicing notation.

[Leave feedback on this section »](#)

9.2 Lists Are Mutable Sequences

The `list` data structure is another sequence type in Python. Just like strings and tuples, lists contain items that are indexed by integers, starting with 0.

On the surface, lists look and behave a lot like tuples. You can use index and slicing notation with lists, check for the existence of an element using `in`, and iterate over lists with a `for` loop.

Unlike tuples, however, lists are **mutable**, meaning you can change

the value at an index even after the list has been created.

In this section, you will learn how to create lists and compare them with tuples.

Creating Lists

A **list literal** looks almost exactly like a tuple literal, except that it is surrounded with square brackets ([and]) instead of parentheses:

```
>>> colors = ["red", "yellow", "green", "blue"]
>>> type(colors)
<class 'list'>
```

When you inspect a list, Python displays it as a list literal:

```
>>> colors
['red', 'yellow', 'green', 'blue']
```

Like tuples, lists values are not required to be of the same type. The list literal `["one", 2, 3.0]` is perfectly valid.

Aside from list literals, you can also use the `list()` built-in to create a new list object from any other sequence. For instance, the tuple `(1, 2, 3)` can be passed to `list()` to create the list `[1, 2, 3]`:

```
>>> list((1, 2, 3))
[1, 2, 3]
```

You can even create a list from a string:

```
>>> list("Python")
['P', 'y', 't', 'h', 'o', 'n']
```

Each letter in the string becomes an element of the list.

There is more useful way to create a list from a string. You can create a list from a string of a comma-separated list of items using the string object's `.split()` method:

```
>>> groceries = "eggs, milk, cheese"
>>> grocery_list = groceries.split(", ")
>>> grocery_list
['eggs', 'milk', 'cheese']
```

The string argument passed to `.split()` is called the separator. By changing the separator you can split strings into lists in numerous ways:

```
>>> # Split string on semi-colons
>>> "a;b;c".split(";");
['a', 'b', 'c']

>>> # Split string on spaces
>>> "The quick brown fox".split(" ")
['The', 'quick', 'brown', 'fox']

>>> # Split string on multiple characters
>>> "abbaabba".split("ba")
['ab', 'ab', '']
```

In the last example above, the string is split around occurrences of the substring "ba", which occurs first at index 2 and again at index 6. The separator has two characters, only the characters at indices 1, 2, 5, and 6 become elements of the list.

`.split()` always returns a string whose length is one more than the number of separators contained in the string. The string "abbaabba" contains two instances of the separator "ba" so the list returned by `split()` has three elements. Since the third separator isn't followed by any other characters, the third element of the list is set to the empty string.

If the separator is not contained in the string at all, `.split()` returns a list with the string as its only element:

```
>>> "abbaabba".split("c")
['abbaabba']
```

In all, you've seen three ways to create a list:

1. A list literal
2. The `list()` built-in
3. The string `.split()` method

Lists support all of the same operations supported by tuples.

Basic List Operations

Indexing and slicing operations work on lists the same way they do on tuples.

You can access list elements using index notation:

```
>>> numbers = [1, 2, 3, 4]
>>> numbers[1]
2
```

You can create a new list from an existing one using slice notation:

```
>>> numbers[1:3]
[2, 3]
```

You can check for the existence of list elements using the `in` operator:

```
>>> # Check existence of an element
>>> "Bob" in numbers
False
```

Because lists are iterable, you can iterate over them with a `for` loop.

```
>>> # Print only the even numbers in the list
>>> for number in numbers:
...     if number % 2 == 0:
```

```
...     print(number)
...
2
4
```

The major difference between lists and tuples is that elements of lists may be changed, but elements of tuples can not.

Changing Elements in a List

Think of a list as a sequence of numbered slots. Each slot holds a value, and every slot must be filled at all times, but you can swap out the value in a given slot with a new one whenever you want.

The ability to swap values in a list for other values is called **mutability**. Lists are **mutable**. The elements of tuples may not be swapped for new values, so tuples are said to be **immutable**.

To swap a value in a list with another, assign the new value to a slot using index notation:

```
>>> colors = ["red", "yellow", "green", "blue"]
>>> colors[0] = "burgundy"
```

The value at index 0 changes from "red" to "burgundy":

```
>>> colors
['burgundy', 'yellow', 'green', 'blue']
```

You can change several values in a list at once with a **slice assignment**:

```
>>> colors[1:3] = ["orange", "magenta"]
>>> colors
['burgundy', 'orange', 'magenta', 'blue']
```

`colors[1:3]` selects the slots with indices 1 and 2. The values in these slots are assigned to "orange" and "magenta", respectively.

The list assigned to a slice does not need to have the same length as the slice. For instance, you can assign a list of three elements to a slice with two elements:

```
>>> colors = ["red", "yellow", "green", "blue"]
>>> colors[1:3] = ["orange", "magenta", "aqua"]
>>> colors
['red', 'orange', 'magenta', 'aqua', 'blue']
```

The values "orange" and "magenta" replace the original values "yellow" and "green" in `colors` at the indices 1 and 2. Then a new slot is created at index 4 and "blue" is assigned to this index. Finally, "aqua" is assigned to index 3.

When the length of the list being assigned to the slice is less than the length of the slice, the overall length of the original list is reduced:

```
>>> colors
['red', 'orange', 'magenta', 'aqua', 'blue']
>>> colors[1:4] = ["yellow", "green"]
>>> colors
['red', 'yellow', 'green', 'blue']
```

The values "yellow" and "green" replace the values "orange" and "magenta" in `colors` at the indices 1 and 2. Then the value at index 3 is replaced with the value "blue". Finally, the slot at index 4 is removed from `colors` entirely.

The above examples show how to change, or **mutate**, lists using index and slice notation. There are also several list methods that you can use to mutate a list.

List Methods For Adding and Removing Elements

Although you can add and remove elements with slice notation, list methods provide a more natural and readable way to mutate a list.

We'll look at several list methods, starting with how to insert a single

value into a list at a specified index.

`list.insert()`

The `list.insert()` method is used to insert a single new value into a list. It takes two parameters, an index *i* and a value *x*, and inserts the value *x* at index *i* in the list.

```
>>> colors = ["red", "yellow", "green", "blue"]
>>> # Insert "orange" into the second position
>>> colors.insert(1, "orange")
>>> colors
['red', 'orange', 'yellow', 'green', 'blue']
```

There are a couple of important observations to make about this example.

The first observation applies to all list methods. To use them, you first write the name of the list you want to manipulate, followed by a dot (.) and then the name of the list method.

So, to use `insert()` on the `colors` list, you must write `colors.insert()`. This works just like string and number methods do.

Next, notice that when the value "orange" is inserted at the index 1, the value "yellow" and all following values are shifted to the right.

If the value for the index parameter of `.insert()` is larger than the greatest index in the list, the value is inserted at the end of the list:

```
>>> colors.insert(10, "violet")
>>> colors
['red', 'orange', 'yellow', 'green', 'blue', 'violet']
```

Here the value "violet" is actually inserted at index 5, even though `.insert()` was called with 10 for the index.

You can also use negative indices with `.insert()`:

```
>>> colors.insert(-1, "indigo")
>>> colors
['red', 'orange', 'yellow', 'green', 'blue', 'indigo', 'violet']
```

This inserts "indigo" into the slot at index -1 which is the last element of the list. The value "violet" is shifted to the right by one slot.

Important

When you `.insert()` an item into a list, you do not need to assign the result to the original list.

For example, the following code actually erases the colors list:

```
>>> colors = colors.insert(-1, "indigo")
>>> print(colors)
None
```

`.insert()` is said to alter `colors` **in place**. This is true for all list methods that do not return a value.

If you can insert a value at a specified index, it only makes sense that you can also remove an element at a specified index.

`list.pop()`

The `list.pop()` method takes one parameter, an index `i`, and removes the value from the list at that index. The value that is removed is returned by the method:

```
>>> color = colors.pop(3)
>>> color
'green'
>>> colors
['red', 'orange', 'yellow', 'blue', 'indigo', 'violet']
```

Here, the value "green" at index 3 is removed and assigned to the variable `color`. When you inspect the `colors` list, you can see that the string "green" has indeed been removed.

Unlike `.insert()`, Python raises an `IndexError` if you pass to `.pop()` an argument larger than the last index:

```
>>> colors.pop(10)
Traceback (most recent call last):
  File "<pyshell#0>", line 1, in <module>
    colors.pop(10)
IndexError: pop index out of range
```

Negative indices also work with `.pop()`:

```
>>> colors.pop(-1)
'violet'
>>> colors
['red', 'orange', 'yellow', 'blue', 'indigo']
```

If you do not pass a value to `.pop()`, it removes the last item in the list:

```
>>> colors.pop()
'indigo'
>>> colors
['red', 'orange', 'yellow', 'blue']
```

This way of removing the final element, by calling `.pop()` with no specified index, is generally considered the most Pythonic.

`list.append()`

The `list.append()` method is used to append a new element to the end of a list:

```
>>> colors.append("indigo")
>>> colors
['red', 'orange', 'yellow', 'blue', 'indigo']
```

After calling `.append()`, the length of the list increases by one and the value "indigo" is inserted into the final slot. Note that `.append()` alters the list in place, just like `.insert()`.

.append() is equivalent to inserting an element at an index greater than or equal to the length of the list. The above example could also have been written as follows:

```
>>> colors.insert(len(colors), "indigo")
```

.append() is both shorter and more descriptive than using .insert() this way, and is generally considered the more Pythonic way of adding an element to the end of a list.

list.extend()

The list.extend() method is used to add several new elements to the end of a list:

```
>>> colors.extend(["violet", "ultraviolet"])
>>> colors
['red', 'orange', 'yellow', 'blue', 'indigo', 'violet', 'ultraviolet']
```

.extend() takes a single parameter that must be an iterable type. The elements of the iterable are appended to the list in the same order that they appear in the argument passed to .extend().

Just like .insert() and .append(), .extend() alters the list in place.

Typically, the argument passed to .extend() is another list, but it could also be a tuple. For example, the above example could be written as follows:

```
>>> colors.extend(("violet", "ultraviolet"))
```

The four list methods discussed in this section make up the most common methods used with lists. The following table serves to recap everything you have seen here:

List Method	Description
.insert(<i>i</i> , <i>x</i>)	Insert the value <i>x</i> at index <i>i</i>
.append(<i>x</i>)	Insert the value <i>x</i> at the end of the list

List Method	Description
.extend(iterable)	Insert all the values of iterable at the end of the list, in order
.pop(i)	Remove and return the element at index i

In addition to list methods, Python has a couple of useful built-in functions for working with lists of numbers.

Lists of Numbers

One very common operation with lists of numbers is to add up all the values to get the total.

You can do this with a `for` loop:

```
>>> nums = [1, 2, 3, 4, 5]
>>> total = 0
>>> for number in nums:
...     total = total + number
...
>>> total
15
```

First you initialize the variable `total` to 0, and then loop over each number in `nums` and add it to `total`, finally arriving at the value 15.

Although this `for` loop is straightforward, there is a much more succinct way of doing this in Python:

```
>>> sum([1, 2, 3, 4, 5])
15
```

The built-in `sum()` function takes a list as an argument and returns the total of all the values in the list.

If the list passed to `sum()` contains any values that aren't numeric, a `TypeError` is raised:

```
>>> sum([1, 2, 3, "four", 5])
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
TypeError: unsupported operand type(s) for +: 'int' and 'str'
```

Besides `sum()`, there are two other useful built-in functions for working with lists of numbers: `min()` and `max()`. These functions return the minimum and maximum values in the list, respectively:

```
>>> min([1, 2, 3, 4, 5])
1

>>> max([1, 2, 3, 4, 5])
5
```

Note that `sum()`, `min()`, and `max()` also work with tuples:

```
>>> sum((1, 2, 3, 4, 5))
15

>>> min((1, 2, 3, 4, 5))
1

>>> max((1, 2, 3, 4, 5))
5
```

The fact that `sum()`, `min()`, and `max()` are all built-in to Python tells you that they are used frequently. Chances are, you'll find yourself using them quite a bit in your own programs!

List Comprehensions

Yet another way to create a list from an existing iterable is with a **list comprehension**:

```
>>> numbers = (1, 2, 3, 4, 5)
>>> squares = [num**2 for num in numbers]
```

```
>>> squares  
[1, 4, 9, 16, 25]
```

A list comprehension is a short-hand for a `for` loop. In the example above, a tuple literal containing five numbers is created and assigned to the `numbers` variable. On the second line, a list comprehension loops over each number in `numbers`, squares each number, and adds it to a new list called `squares`.

To create the `sqaures` list using a traditional `for` loop involves first creating an empty list, looping over the numbers in `numbers`, and appending the square of each number to the list:

```
>>> squares = []  
>>> for num in numbers:  
...     squares.append(num**2)  
...  
>>> squares  
[1, 4, 9, 16, 25]
```

List comprehensions are commonly used to convert values in one list to a different type.

For instance, suppose you needed to convert a list of strings containing floating point values to a list of `float` objects. The following list comprehensions achieves this:

```
>>> str_numbers = ["1.5", "2.3", "5.25"]  
>>> float_numbers = [float(value) for value in str_numbers]  
>>> float_numbers  
[1.5, 2.3, 5.25]
```

List comprehensions are not unique to Python, but they are one of its many beloved features. If you find yourself creating an empty list, looping over some other iterable, and appending new items to the list, then chances are you can replace your code with a list comprehension!

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Create a list named `food` with two elements "rice" and "beans".
2. Append the string "broccoli" to `food` using `.append()`.
3. Add the string "bread" and "pizza" to "food" using `.extend()`.
4. Print the first two items in the `food` list using `print()` and slicing notation.
5. Print the last item in `food` using `print()` and index notation.
6. Create a list called `breakfast` from the string "eggs, fruit, orange juice" using the string `.split()` method.
7. Verify that `breakfast` has three items using `len()`.
8. Create a new list called `lengths` using a list comprehension that contains the lengths of each string in the `breakfast` list.

[Leave feedback on this section »](#)

9.3 Nesting, Copying, and Sorting Tuples and Lists

Now that you have learned what tuples and lists are, how to create them, and some basic operations with them, let's look at three more concepts:

1. Nesting
2. Copying
3. Sorting

Nesting Lists and Tuples

Lists and tuples can contain values of any type. That means lists and tuples can contain lists and tuples as values. A **nested list**, or **nested**

tuple, is a list or tuple that is contained as a value in another list or tuple.

For example, the following list has two values, both of which are other lists:

```
>>> two_by_two = [[1, 2], [3, 4]]  
  
>>> # two_by_two has length 2  
>>> len(two_by_two)  
2  
  
>>> # Both elements of two_by_two are lists  
>>> two_by_two[0]  
[1, 2]  
>>> two_by_two[1]  
[3, 4]
```

Since `two_by_two[1]` returns the list `[3, 4]`, you can use **double index notation** to access an element in the nested list:

```
>>> two_by_two[1][0]  
3
```

First, Python evaluates `two_by_two[1]` and returns `[3, 4]`. Then Python evaluates `[3, 4][0]` and returns the first element 3.

In very loose terms, you can think of a list of lists or a tuple of tuples as a sort of table with rows and columns.

The `two_by_two` list has two rows, `[1, 2]` and `[3, 4]`. The columns are made of the corresponding elements of each row, so the first column contains the elements 1 and 3, and the second column contains the elements 2 and 4.

This table analogy is only an informal way of thinking about a list of lists, though. For example, there is no requirement that all the lists in a list of lists have the same length, in which case this table analogy starts to break down.

Note

Readers interested in data analysis or scientific computing may recognize lists of lists as a sort of matrix of values.

While you can use the built in `list` and `tuple` types for matrices, better alternatives exist. To learn how to work with matrices in Python, check out Chapter 17.

Copying a List

Sometimes you need to copy one list into another list. However, you can't just reassign one list object to another list object, because you'll get this (possibly surprising) result:

```
>>> animals = ["lion", "tiger", "frumious Bandersnatch"]
>>> large_cats = animals
>>> large_cats.append("Tigger")
>>> animals
['lion', 'tiger', 'frumious Bandersnatch', 'Tigger']
```

In this example, you first assign the list stored in the `animals` variable to the variable `large_cats`, and then we add a new string to the `large_cats` list. But, when the contents of `animals` are displayed you can see that the original list has also been changed.

This is a quirk of object-oriented programming, but it's by design. When you say `large_cats = animals`, the `large_cats` and `animals` variables both refer to the same object.

A variable name is really just a reference to a specific location in computer memory. Instead of copying all the contents of the list object and creating a new list, `large_cats = animals` assigns the memory location referenced by `animals` to `large_cats`. That is, both variables now refer to the same object in memory, and any changes made to one will affect the other.

To get an independent copy of the `animals` list, you can use slicing no-

tation to return a new list with the same values:

```
>>> animals = ["lion", "tiger", "frumious Bandersnatch"]
>>> large_cats = animals[:]
>>> large_cats.append("leopard")
>>> large_cats
['lion', 'tiger', 'frumious Bandersnatch', 'leopard']
>>> animals
["lion", "tiger", "frumious Bandersnatch"]
```

Since no index numbers are specified in the `[:]` slice, every element of the list is returned from beginning to end. The `large_cats` list now has the same elements as `animals`, and in the same order, but you can `.append()` items to it without changing the list assigned to `animals`.

If you want to make a copy of a list of lists, you can do so using the `[:]` notation you saw earlier:

```
>>> matrix1 = [[1, 2], [3, 4]]
>>> matrix2 = matrix1[:]
>>> matrix2[0] = [5, 6]
>>> matrix2
[[5, 6], [3, 4]]
>>> matrix1
[[1, 2], [3, 4]]
```

Let's see what happens when you change the first element of the second list in `matrix2`:

```
>>> matrix2[1][0] = 1
>>> matrix2
[[5, 6], [1, 4]]
>>> matrix1
[[1, 2], [1, 4]]
```

Notice that the second list in `matrix1` was also altered!

This happens because a list does not really contain objects themselves, but references to those objects in memory. When you make a copy

of the list using the `[:]` notation, a new list is returned containing the same references as the original list. In programming jargon, this method of copying a list is called a **shallow copy**.

To make a copy of both the list and all of the elements it contains, you must use what is known as a **deep copy**. This method of copying is beyond the scope of this course. For more information on shallow and deep copies, check out the [Shallow vs Deep Copying of Python Objects](#) article on [realpython.com](#).

Sorting Lists

Lists have a `.sort()` method that sorts all of the items in ascending order. By default, the list is sorted in alphabetical or numerical order, depending on the type of elements in the list:

```
>>> # Lists of strings are sorted alphabetically
>>> colors = ["red", "yellow", "green", "blue"]
>>> colors.sort()
>>> colors
['blue', 'green', 'red', 'yellow']

>>> # Lists of numbers are sorted numerically
>>> numbers = [1, 10, 5, 3]
>>> numbers.sort()
>>> numbers
[1, 3, 5, 10]
```

Notice that `.sort()` sorts the list in place, so you don't need to assign its result to anything.

`.sort()` has an option parameter called `key` that can be used to adjust how the list gets sorted. The `key` parameter accepts a function, and the list is sorted based on the return value of that function.

For example, to sort a list of strings by the length of each string, you can pass the `len` function to `key`:

```
>>> colors = ["red", "yellow", "green", "blue"]
>>> colors.sort(key=len)
>>> colors
['red', 'blue', 'green', 'yellow']
```

You don't need to call the function when you pass it to key. Pass the name of the function without any parentheses. For instance, in the previous example the name `len` is passed to key, and not `len()`.

The function that gets passed to key must only accept a single argument.

You can also pass user defined functions to key. In the following example, a function called `get_second_element()` is used to sort a list of tuples by their second elements:

```
>>> def get_second_element(item):
...     return item[1]
...
>>> items = [(4, 1), (1, 2), (-9, 0)]
>>> items.sort(key=get_second_element)
>>> items
[(-9, 0), (4, 1), (1, 2)]
```

Keep in mind that any function that you pass to key must accept only a single argument.

Review Exercises

1. Create a tuple `data` with two values. The first value should be the tuple `(1, 2)` and the second value should be the tuple `(3, 4)`.
2. Write a for loop that loops over `data` and prints the sum of each nested tuple. The output should look like this:

```
Row 1 sum: 3
Row 2 sum: 7
```

3. Create the following list `[4, 3, 2, 1]` and assign it to the variable

numbers.

4. Create a copy of the numbers list using the [:] slicing notation.
5. Sort the numbers list in numerical order using the .sort() method.

[Leave feedback on this section »](#)

9.4 Challenge: List of lists

Write a program that contains the following lists of lists:

```
universities = [  
    ['California Institute of Technology', 2175, 37704],  
    ['Harvard', 19627, 39849],  
    ['Massachusetts Institute of Technology', 10566, 40732],  
    ['Princeton', 7802, 37000],  
    ['Rice', 5879, 35551],  
    ['Stanford', 19535, 40569],  
    ['Yale', 11701, 40500]  
]
```

Define a function, enrollment_stats(), that takes, as an input, a list of lists where each individual list contains three elements: (a) the name of a university, (b) the total number of enrolled students, and (c) the annual tuition fees.

enrollment_stats() should return two lists: the first containing all of the student enrollment values and the second containing all of the tuition fees.

Next, define a mean() and a median() function. Both functions should take a single list as an argument and return the mean and median of the values in each list.

Using universities, enrollment_stats(), mean(), and median(), calculate the total number of students, the total tuition, the mean and median of the number of students, and the mean and median tuition values.

Finally, output all values, and format the output so that it looks like this:

```
*****
Total students: 77,285
Total tuition: $ 271,905

Student mean:    11,040.71
Student median: 10,566

Tuition mean:   $ 38,843.57
Tuition median: $ 39,849
*****
```

You can find the solutions to this code challenge and many other bonus resources online at realpython.com/python-basics/resources.

[Leave feedback on this section »](#)

9.5 Challenge: Wax Poetic

In this challenge, you'll write a program that generates poetry.

Create five lists for different word types:

- Nouns: `["fossil", "horse", "aardvark", "judge", "chef", "mango", "extrovert", "gorilla"]`
- Verbs: `["kicks", "jingles", "bounces", "slurps", "meows", "explodes", "curdles"]`
- Adjectives: `["furry", "balding", "incredulous", "fragrant", "exuberant", "glistening"]`
- Prepositions: `["against", "after", "into", "beneath", "upon", "for", "in", "like", "over", "within"]`
- Adverbs: `["curiously", "extravagantly", "tantalizingly", "furiously", "sensuously"]`

Randomly select the following number of elements from each list:

- 3 nouns
- 3 verbs
- 3 adjectives
- 2 prepositions
- 1 adverb

You can do this with the `choice()` function in the `random` module. This function takes a list as input and returns a randomly selected element of the list.

For example, here's how you use `random.choice()` to get random element from the list `["a", "b", "c"]`:

```
import random

random_element = random.choice(["a", "b", "c"])
```

Using the randomly selected words, generate and display a poem with the following structure inspired by [Clifford Pickover](#):

```
{A/An} {adj1} {noun1}

{A/An} {adj1} {noun1} {verb1} {prep1} the {adj2} {noun2}
{adverb1}, the {noun1} {verb2}
the {noun2} {verb3} {prep2} a {adj3} {noun3}
```

Here, `adj` stands for adjective and `prep` for preposition.

Here's an example of the kind of poem your program might generate:

```
A furry horse

A furry horse curdles within the fragrant mango
extravagantly, the horse slurps
the mango meows beneath a balding extrovert
```

Every time your program runs, it should generate a new poem.

You can find the solutions to this code challenge and many other bonus resources online at realpython.com/python-basics/resources.

[Leave feedback on this section »](#)

9.6 Store Relationships in Dictionaries

One of the most useful data structures in Python is the **dictionary**.

In this section, you'll learn what a dictionary is, how dictionaries differ from lists and tuples, and how to define and use dictionaries in your own code.

What is a Dictionary?

In plain English, a dictionary is a book containing the definitions of words. Each entry in a dictionary has two parts: the word being defined, and its definition.

Python dictionaries, like lists and tuples, store a collection of objects. However, instead of storing objects in a sequence, dictionaries hold information in pairs of data called **key-value pairs**. That is, each object in a dictionary has two parts: a **key** and a **value**.

The **key** in a key-value pair is a unique name that identifies the **value** part of the pair. Comparing this to an English dictionary, the key is like the word being defined and the value is like the definition of the word.

For example, you could use a dictionary to store names of states and their capitals:

Key	Value
"California"	"Sacramento"
"New York"	"Albany"

Key	Value
"Texas"	"Austin"

In the table above, the keys of the dictionary are the names of the states, and the values of the dictionary are the names of the capitals.

The difference between an English dictionary and a Python dictionary is that the relationship between a key and its value is completely arbitrary. Any key can be assigned to any value.

For example, the following table of key-value pairs is valid:

Key	Value
1	"Sunday"
"red"	12:45pm
17	True

The keys in this table don't appear to be related to the values at all. The only relationship is that each key is assigned to its corresponding value by the dictionary.

In this sense, a Python dictionary is much more like a **map** than it is an English dictionary. The term map here comes from mathematics. It is used to describe a relation between two sets of values, not a geographical map.

In practice, it is this idea of dictionaries as a map that is particularly useful. Under this lens, the English dictionary is a special case of a map that relates words to their definitions.

So in summary, a Python dictionary is a data structure that relates a set of keys to a set of values. Each key is assigned a single value, which defines a relationship between the two sets.

Now that you have an idea what a dictionary is, let's see how to create

dictionaries in Python code.

Creating Dictionaries

The following code creates a **dictionary literal** containing names of states and their capitals:

```
>>> capitals = {  
    "California": "Sacramento",  
    "New York": "Albany",  
    "Texas": "Austin",  
}
```

Notice that each key is separated from its value by a colon (:), each key-value pair is separated by a comma (,), and the entire dictionary is enclosed in curly braces ({ and }).

You can also create a dictionary from a sequence of tuples using the `dict()` built-in:

```
>>> key_value_pairs = (  
...     ("California", "Sacramento"),  
...     ("New York", "Albany"),  
...     ("Texas", "Austin"),  
)  
>>> capitals = dict(key_value_pairs)
```

When you inspect a dictionary, it is displayed as a dictionary literal, regardless of how it was created:

```
>>> capitals  
{'California': 'Sacramento', 'New York': 'Albany', 'Texas': 'Austin'}
```

Note

If you happen to be following along with a Python version older than 3.6, then you will notice that the output dictionaries in the interactive window have a different order than the ones that appear in these examples.

Prior to Python 3.6, the order of key-value pairs in a Python dictionary was random. In later versions, the order of the key-value pairs is guaranteed to match the order in which they were inserted.

You can create an empty dictionary using either a literal or `dict()`:

```
>>> {}
{}
>>> dict()
{}
```

Now that we've created a dictionary, let's look at how you access its values.

Accessing Dictionary Values

To access a value in a dictionary, enclose the corresponding key in square brackets ([and]) at the end of dictionary or a variable name assigned to a dictionary:

```
>>> capitals["Texas"]
'Austin'
```

The bracket notation used to access a dictionary value looks similar to the index notation used to get values from strings, lists, and tuples. However, dictionaries are a fundamentally different data structure than sequence types like lists and tuples.

To see the difference, let's step back for a second and notice that we could just as well define the `capitals` dictionary as a list:

```
>>> capitals_list = ["Sacramento", "Albany", "Austin"]
```

You can use index notation to get the capital of each of the three states from the `capitals` dictionary:

```
>>> capitals_list[0] # Capital of California  
'Sacramento'  
  
>>> capitals_list[2] # Capital of Texas  
'Austin'
```

One nice thing about dictionaries is that they can be used to provide context to the values they contain. Typing `capitals["Texas"]` is easier to understand than `capitals_list[2]`, and you don't have to remember the order of data in a long list or tuple.

This idea of ordering is really the main difference between how items in a sequence type are accessed compared to a dictionary.

Values in a sequence type are accessed by index, which is an integer value expressing the order of items in the sequence.

On the other hand, items in a dictionary are accessed by a key, which doesn't define any kind of order, but just provides a label that can be used to reference the value.

Adding and Removing Values in a Dictionary

Like lists, dictionaries are mutable data structures. This means you can add and remove items from a dictionary.

Let's add the capital of Colorado to the `capitals` dictionary:

```
>>> capitals["Colorado"] = "Denver"
```

First you use the square bracket notation with "Colorado" as the key, as if you were looking up the value. Then you use the assignment operator `=` to assign the value "Denver" to the new key.

When you inspect `capitals`, you see that a new key "Colorado" exists with the value "Denver":

```
>>> capitals
{'California': 'Sacramento', 'New York': 'Albany', 'Texas': 'Austin',
'Colorado': 'Denver'}
```

Each key in a dictionary can only be assigned a single value. If a key is given a new value, Python just overwrites the old one:

```
>>> capitals["Texas"] = "Houston"
>>> capitals
{'California': 'Sacramento', 'New York': 'Albany', 'Texas': 'Houston',
'Colorado': 'Denver'}
```

To remove an item from a dictionary, use the `del` keyword with the key for the value you want to delete:

```
>>> del capitals["Texas"]
>>> capitals
{'California': 'Sacramento', 'New York': 'Albany',
'Colorado': 'Denver'}
```

Checking the Existence of Dictionary Keys

If you try to access a value in a dictionary using a key that doesn't exist, Python raises a `KeyError`:

```
>>> capitals["Arizona"]
Traceback (most recent call last):
  File "<pyshell#1>", line 1, in <module>
    capitals["Arizona"]
KeyError: 'Arizona'
```

The `KeyError` is the most common error encountered when working with dictionaries. Whenever you see it, it means that an attempt was made to access a value using a key that doesn't exist.

You can check that a key exists in a dictionary using the `in` keyword:

```
>>> "Arizona" in capitals  
False  
>>> "California" in capitals  
True
```

With `in`, you can first check that a key exists before doing something with the value for that key:

```
>>> if "Arizona" in capitals:  
...     # Only print if the "Arizona" key exists  
...     print(f"The capital of Arizona is {capitals['Arizona']}")
```

It is important to remember that `in` only checks the existence of keys:

```
>>> "Sacramento" in capitals  
False
```

Even though "Sacramento" is a value for the existing "California" key in `capitals`, checking for its existence returns `False`.

Iterating Over Dictionaries

Like lists and tuples, dictionaries are iterable. However, looping over a dictionary is a bit different than looping over a list or tuple.

When you loop over a dictionary with a `for` loop, you iterate over the dictionary's keys:

```
>>> for key in capitals:  
...     print(key)  
...  
California  
New York  
Colorado
```

So, if you want to loop over the `capitals` dictionary and print “The capital of X is Y”, where X is the name of the state and Y is the state's capital, you can do the following:

```
>>> for state in capitals:  
    print(f"The capital of {state} is {capitals[state]}")
```

```
The capital of California is Sacramento  
The capital of New York is Albany  
The capital of Colorado is Denver
```

However, there is a slightly more succinct way to do this using the `.items()` dictionary method. `.items()` returns a list-like object containing tuples of key-value pairs. For example, `capitals.items()` returns a list of tuples of states and their corresponding capitals:

```
>>> capitals.items()  
dict_items([('California', 'Sacramento'), ('New York', 'Albany'),  
('Colorado', 'Denver')])
```

The object returned by `.items()` isn't really a list. It has a special type called a `dict_items`:

```
>>> type(capitals.items())  
<class 'dict_items'>
```

You don't need to worry about what `dict_items` really is, because you usually won't work with it directly. The important thing to know about it is that you can use `.items()` to loop over a dictionary's keys and values simultaneously.

Let's rewrite the previous loop using `.items()`:

```
>>> for state, capital in capitals.items():  
...     print(f"The capital of {state} is {capital}")
```

```
The capital of California is Sacramento  
The capital of New York is Albany  
The capital of Colorado is Denver
```

When you loop over `capitals.items()`, each iteration of the loop produces a tuple containing the state name and the corresponding capital

city name. By assigning this tuple to `state`, `capital`, the components of the tuple are unpacked into the two variable `state` and `capital`.

Dictionary Keys and Immutability

In the `capitals` dictionary you've been working with throughout this section, each key is a string. However, there is no rule that says dictionary keys must all be of the same type.

For instance, you can add an integer key to `capitals`:

```
>>> capitals[50] = "Honolulu"
>>> capitals
{'California': 'Sacramento', 'New York': 'Albany',
 'Colorado': 'Denver', 50: 'Honolulu'}
```

There is only one restriction on what constitutes a valid dictionary key. Only immutable types are allowed. This means, for example, that a list cannot be a dictionary key.

Consider this: what should happen if a list were used as a key in a dictionary and, somewhere later in the code, the list is changed?

Should the list be associated to the same value as the old list in the dictionary? Or should the value for the old key be removed from the dictionary all together?

Rather than make a guess about what should be done, Python raises an exception:

```
>>> capitals[[1, 2, 3]] = "Bad"
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
TypeError: unhashable type: 'list'
```

It might not seem fair that some types can be keys and others can't, but it's important that a programming language always has well-defined behavior. It should never make guesses about what the author intended.

For reference, here's a list of all the data types you've learned about so far that are valid dictionary keys:

Valid Dictionary Key Types

integers
floats
strings
booleans
tuples

Unlike keys, dictionary values can be any valid Python type, including other dictionaries!

Nested Dictionaries

Just as you can nest lists inside of other lists, and tuples inside of other tuples, you can create nested dictionaries.

Let's alter the `capitals` dictionary to illustrate this idea. Instead of mapping state names to their capital cities, we'll create a dictionary that maps each state name to a dictionary containing the capital city and the state flower.

```
>>> states = {
...     "California": {
...         "capital": "Sacramento",
...         "flower": "California Poppy"
...     },
...     "New York": {
...         "capital": "Albany",
...         "flower": "Rose"
...     },
...     "Texas": {
...         "capital": "Austin",
...         "flower": "Bluebonnet"
...     }
... }
```

```
...     },
... }
```

The value of each key is a dictionary:

```
>>> states["Texas"]
{'capital': 'Austin', 'flower': 'Bluebonnet'}
```

To get the Texas state flower, first get the value at the key "Texas", and then the value at the key "flower":

```
>>> states["Texas"]["flower"]
'Bluebonnet'
```

Nested dictionaries come up more often than you might expect. They are particularly useful when working with data transmitted over the web. Nested dictionaries are also great for modeling structured data, such as spreadsheets or relational databases.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Create an empty dictionary named `captains`.
2. Using the square bracket notation, enter the following data into the dictionary, one item at a time:

```
'Enterprise': 'Picard'
'Voyager': 'Janeway'
'Defiant': 'Sisko'
```
3. Write two `if` statements that check if "Enterprise" and "Discovery" exist as keys in the dictionary. Set their values to "unknown" if the key does not exist.
4. Write a `for` loop to display the ship and captain names contained in the dictionary. For example, the output should look something like this:

The Enterprise is captained by Picard.

5. Delete "Discovery" from the dictionary.
6. Bonus: Make the same dictionary by using `dict()` and passing in the initial values when you first create the dictionary.

[Leave feedback on this section »](#)

9.7 Challenge: Capital City Loop

Review your state capitals along with dictionaries and `while` loops!

First, finish filling out the following dictionary with the remaining states and their associated capitals in a file called `capitals.py`.

```
capitals_dict = {  
    'Alabama': 'Montgomery',  
    'Alaska': 'Juneau',  
    'Arizona': 'Phoenix',  
    'Arkansas': 'Little Rock',  
    'California': 'Sacramento',  
    'Colorado': 'Denver',  
    'Connecticut': 'Hartford',  
    'Delaware': 'Dover',  
    'Florida': 'Tallahassee',  
    'Georgia': 'Atlanta',  
}
```

Next, pick a random state name from the dictionary, and assign both the state and it's capital to two variables. You'll need to `import` the `random` module at the top of your program.

Then display the name of the state to the user and ask them to enter the capital. If the user answers, incorrectly, repeatedly ask them for the capital name until they either enter the correct answer or type the word "exit".

If the user answers correctly, display "Correct" and end the program.

However, if the user exits without guessing correctly, display the correct answer and the word "Goodbye".

Note

Make sure the user is not punished for case sensitivity. In other words, a guess of "Denver" is the same as "denver". Do the same for exiting—"EXIT" and "Exit" should work the same as "exit".

You can find the solutions to this code challenge and many other bonus resources online at realpython.com/python-basics/resources.

[Leave feedback on this section »](#)

9.8 How to Pick a Data Structure

In this chapter, you've learned about three data structures native to Python: lists, tuples, and dictionaries.

You might be wondering, "How do I know when to use which data structure?" It's a great question, and one many new Python programmers struggle with.

The type of data structure you use depends on the problem you are solving, and there is no hard and fast rule you can use to pick the right data structure every time. You'll always need to spend a little time thinking about the problem, and which structure works best for it.

Fortunately, there are some guidelines you can use to help you make the right choice. These are presented below:

Use a list when:

- Data has a natural order to it
- You will need to update or alter the data during the program
- The primary purpose of the data structure is iteration

Use a tuple when:

- Data has a natural order to it
- You **will not** need to update or alter the data during the program
- The primary purpose of the data structure is iteration

Use a dictionary when:

- The data is unordered, or the order does not matter
- You will need to update or alter the data during the program
- The primary purpose of the data structure is looking up values

[Leave feedback on this section »](#)

9.9 Challenge: Cats With Hats

You have 100 cats.

One day you decide to arrange all your cats in a giant circle. Initially, none of your cats have any hats on. You walk around the circle 100 times, always starting at the same spot, with the first cat (cat # 1). Every time you stop at a cat, you either put a hat on it if it doesn't have one on, or you take its hat off if it has one on.

1. The first round, you stop at every cat, placing a hat on each one.
2. The second round, you only stop at every second cat (#2, #4, #6, #8, etc.).
3. The third round, you only stop at every third cat (#3, #6, #9, #12, etc.).
4. You continue this process until you've made 100 rounds around the cats (e.g., you only visit the 100th cat).

Write a program that simply outputs which cats have hats at the end.

Note

This is not an easy problem by any means. Honestly, the code is simple. This problem is often seen on job interviews as it tests your ability to reason your way through a difficult problem. Stay calm. Start with a diagram, and then write pseudo code. Find a pattern. Then code!

You can find the solutions to this code challenge and many other bonus resources online at realpython.com/python-basics/resources.

[Leave feedback on this section »](#)

9.10 Summary and Additional Resources

In this chapter, you learned about three data structures: lists, tuples, and dictionaries.

Lists, such as [1, 2, 3, 4], are mutable sequences of objects. You can interact with lists using various list methods, such as .append(), .remove(), and .extend(). Lists can be sorted using the .sort() method. You can access individual elements of a list using subscript notation, just like strings. Slicing notation also works with lists.

Tuples, like lists, are sequences of objects. The big difference between lists and tuples is that tuples are immutable. Once you create a tuple, it cannot be changed. Just like lists, you can access elements by index and using slicing notation.

Dictionaries store data as key-value pairs. They are not sequences, so you cannot access elements by index. Instead, you access elements by their key. Dictionaries are great for storing relationships, or when you need quick access to data. Like lists, dictionaries are mutable.

Lists, tuples and dictionaries are all iterable, meaning they can be looped over. You saw how to loop over all three of these structures

using for loops.

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-9

Additional Resources

To learn more about lists, tuples, and dictionaries, check out the following resources:

- [Lists and Tuples in Python](#)
- [Dictionaries in Python](#)
- Recommended resources on [realpython.com](#)

[Leave feedback on this section »](#)

Chapter 10

Object-Oriented Programming (OOP)

OOP, or **Object-Oriented Programming**, is a method of structuring a program by bundling related properties and behaviors into individual **objects**.

Conceptually, objects are like components of a system. Think of a program as a factory assembly line of sorts. A system component at each step of the assembly line processes some material a little bit, ultimately transforming raw material into a finished product.

An object contains data, like the raw or pre-processed materials at each step on an assembly line, and behavior, like the action each assembly line component performs.

In this chapter, you will learn how to:

- Create a `class`, which is like a blueprint for creating an object
- Use classes to create new objects
- Model systems with class inheritance

Let's get started!

[Leave feedback on this section »](#)

10.1 Define a Class

Primitive data structures—like numbers, strings, and lists—are designed to represent simple things, such as the cost of something, the name of a poem, and your favorite colors, respectively. What if you want to represent something much more complicated?

For example, let's say you wanted to track employees in an organization. You need to store some basic information about each employee, such as their name, age, position, and the year they started working.

One way to do this is to represent each employee as a list:

```
kirk = ["James Kirk", 34, "Captain", 2265]
spock = ["Spock", 35, "Science Officer", 2254]
mccoy = ["Leonard McCoy", "Chief Medical Officer", 2266]
```

There are a number of issues with this approach.

First, when you reference `kirk[0]` several lines away from where the `kirk` list is declared, will you remember that the 0th element of the list is the employee's name? What if not every employee has the same number of elements in the list?

Second, in the `mccoy` list above, the age is missing, so `mccoy[1]` will return "Chief Medical Officer" instead of Dr. McCoy's age.

A great way to make this type of code more manageable and more maintainable is to use **classes**.

Classes vs Instances

Classes are used to create user-defined data structures. Classes also have special functions, called **methods**, that define behaviors and actions that an object created from the class can perform with its data.

In this chapter you'll create a `Dog` class that stores some basic information about a dog.

It's important to note that a class just provides structure. A class is a blueprint for how something should be defined. It doesn't actually provide any real content itself. The `Dog` class may specify that the name and age are necessary for defining a dog, but it will not actually state what a specific dog's name or age is.

While the class is the blueprint, an **instance** is an object built from a class that contains real data. An instance of the `Dog` class is not a blueprint anymore. It's an actual dog with a name, like Miles, who's four years old.

Put another way, a class is like a form or questionnaire. It defines the needed information. After you fill out the form, your specific copy is an instance of the class. It contains actual information relevant to you.

You can fill out multiple copies of a form to create many different instances, but without the form as a guide, you would be lost, not knowing what information is required. Thus, before you can create individual instances of an object, you must first specify what is needed by defining a class.

How to Define a Class

All class definitions start with the `class` keyword, which is followed by the name of the class and a colon. This is similar to the signature of a function, except that you don't need to add any parameters in parentheses. Any code that is indented below the class definition is considered part of the class's body.

Here is an example of a simple `Dog` class:

```
class Dog:  
    pass
```

The body of the `Dog` class consists of a single statement: the `pass` keyword. `pass` is often used as a place holder where code will eventually go. It allows you to run this code without throwing an error.

Note

Unlike functions and variables, the convention for naming classes in Python is to use **CamelCase notation**, starting with a capital letter. For example, a class for a specific breed of a dog, like the Jack Russell Terrier, would be written as `JackRussellTerrier`.

The `Dog` class isn't very interesting right now, so let's spruce it up a bit by defining some properties that all `Dog` objects should have. There are a number of properties that we can choose from, such as name, age, coat color, and breed. To keep things simple, we'll stick with just two for now: name and age.

To define the properties, or **instance attributes**, that all `Dog` objects must have, you need to define a special method called `__init__()`. This method is run every time a new `Dog` object is created and tells Python what the initial **state**—that is, the initial values of the object's properties—of the object should be.

The first positional argument of `__init__()` is always a variable that references the class instance. This variable is almost universally named `self`. After the `self` argument, you can specify any other arguments required to create an instance of the class.

The following updated definition of the `Dog` class shows how to write an `__init__()` method that creates two instance attributes: `.name` and `.age`:

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

Notice that the function signature—the part that starts with the `def` keyword—is indented four spaces. The body of the function is indented by eight spaces. This indentation is vitally important. It tells Python that the `__init__()` method belongs to the `Dog` class.

Without the indentation, Python would treat `__init__()` as just another function.

Note

Functions that belong to a class are called **instance methods** because they belong to the instance of a class. For example, `list.append()` and `string.find()` are instance methods.

In the body of the `__init__()` method, there are two statements using the `self` variable. The first line, `self.name = name`, creates an instance attribute called `name` and assigns to it the value of the `name` variable that was passed to the `__init__()` method. The second line creates an instance attribute called `age` and assigns to it the value of the `age` argument.

This might look kind of strange. The `self` variable is referring to an instance of the `Dog` class, but we haven't actually created an instance yet. It is a place holder that is used to build the blueprint. Remember, the class is used to define the `Dog` data structure. It does not actually create any instances of individual dogs with specific names and ages.

While instance attributes are specific to each object, **class attributes** are the same for all instances—which in this case is all dogs. In the next example, a class attribute called `species` is created and assigned the value "`Canis familiaris`":

```
class Dog:  
    # Class Attribute  
    species = "Canis familiaris"  
  
    def __init__(self, name, age):  
        self.name = name  
        self.age = age
```

Class attributes are defined directly underneath the first line of the class and outside of any method definition. They must be assigned a value because they are created on a class instance without arguments

to determine what their initial value should be.

You should use class attributes whenever a property should have the same initial value for all instances of a class. Use instance attributes for properties that must be specified before an instance is created.

Now that we have a `Dog` class, let's create some dogs!

[Leave feedback on this section »](#)

10.2 Instantiate an Object

Once a class has been defined, you have a blueprint for creating—also known as **instantiating**—new objects. To instantiate an object, type the name of the class, in the original CamelCase, followed by parentheses containing any values that must be passed to the class's `__init__()` method.

Let's take a look at an actual example. Open IDLE's interactive window and type the following:

```
>>> class Dog:  
...     pass  
...
```

This creates a new `Dog` class with no attributes and methods.

Next, instantiate a new `Dog` object:

```
>>> Dog()  
<__main__.Dog object at 0x106702d30>
```

The output indicates that you now have a new `Dog` object at memory address `0x106702d30`. Note that the address you see on your screen will very likely be different from the address shown here.

Now let's instantiate another `Dog` object:

```
>>> Dog()
<__main__.Dog object at 0x0004ccc90>
```

The new `Dog` instance is located at a different memory address. This is because it is an entirely new instance, completely unique from the first `Dog` object you instantiated.

To see this another way, type the following:

```
>>> a = Dog()
>>> b = Dog()
>>> a == b
False
```

Two new `Dog` objects are created and assigned to the variables `a` and `b`. When `a` and `b` are compared using the `==` operator, the result is `False`. For user defined classes, the default behavior of the `==` operator is to compare the memory addresses of two objects and return `True` if the address is the same and `False` otherwise.

What this means is that even though the `a` and `b` objects are both instances of the `Dog` class and have the exact same attributes and methods—namely, no attributes or methods, in this case—`a` and `b` represent two distinct objects in memory.

Note

The default behavior of the `==` operator can be overridden. How this is done is outside the scope of this book.

If you would like more information on how to customize the behavior of your classes, check out Real Python's [Operator and Function Overloading in Custom Python Classes](#) tutorial.

You can use the `type()` function to determine an object's class:

```
>>> type(a)
<class '__main__.Dog'>
```

Of course, even though both `a` and `b` are distinct `Dog` instances, they have the same type:

```
>>> type(a) == type(b)
```

```
True
```

Class and Instance Attributes

Let's look at a slightly more complex example using the `Dog` class we defined with `.name` and `.age` instance attributes:

```
>>> class Dog:  
...     species = "Canis familiaris"  
...     def __init__(self, name, age):  
...         self.name = name  
...         self.age = age  
...  
>>> buddy = Dog("Buddy", 9)  
>>> miles = Dog("Miles", 4)
```

After declaring the new `Dog` class, two new instances are created—one `Dog` whose name is Buddy and is nine years old, and another named Miles who is four years old.

Does anything look a little strange about how the `Dog` objects are instantiated? The `__init__()` method has three parameters, so why are only two arguments passed to it in the example?

When you instantiate a `Dog` object, Python creates a new instance and passes it to the first parameter of `__init__()`. This essentially removes the `self` parameter, so you only need to worry about the `name` and `age` parameters.

After the `Dog` instances are created, you can access their instance attributes by using **dot notation**:

```
>>> buddy.name
```

```
'Buddy'
```

```
>>> buddy.age
```

```
9
```

```
>>> miles.name
```

```
'Miles'
```

```
>>> miles.age
```

```
4
```

Class attributes are accessed the same way:

```
>>> buddy.species
```

```
'Canis familiaris'
```

One of the biggest advantages of using classes to organize data is that instances are guaranteed to have the attributes you expect:

```
>>> buddy.species == miles.species
```

```
True
```

Both `buddy` and `miles` have the `.species` attribute. Contrast this to the method of using lists to represent similar data structures that you saw at the beginning of the previous section. With a class you no longer have to worry that an attribute may be missing.

Both instance and class attributes can be modified dynamically:

```
>>> buddy.age = 10
```

```
>>> buddy.age
```

```
10
```

```
>>> miles.species = "Felis silvestris"
```

```
>>> miles.species
```

```
'Felis silvestris'
```

In this example, the `.age` attribute of the `buddy` object is changed to 10. Then the `.species` attribute of the `miles` object is changed to "Felis silvestris", which is a species of cat. That makes Miles a pretty strange dog, but it is valid Python!

The important takeaway here is that custom objects are mutable by default. Recall that an object is mutable if it can be altered dynamically. For example, lists and dictionaries are mutable, but strings and tuples are not—they are immutable.

Now that you know the difference between a class and an instance, how to create instances and set class and instance attributes, the next step is to look at instance methods in more detail.

Instance Methods

Instance methods are functions defined inside of a class. This means that they only exist within the context of the object itself and cannot be called without referencing the object. Just like `__init__()`, the first argument of an instance method is always `self`:

```
class Dog:
    species = "Canis familiaris"

    def __init__(self, name, age):
        self.name = name
        self.age = age

        # Instance method
    def description(self):
        return f"{self.name} is {self.age} years old"

        # Another instance method
    def speak(self, sound):
        return f"{self.name} says {sound}"
```

In this example, two new instance methods are defined:

1. `.description()` returns a string displaying the name and age of the dog
2. `.speak()` has one parameter called `sound` and returns a string containing the dog's name and the sound the dog makes.

Let's see how instance methods work in practice. To avoid typing out the whole class in the interactive window, you can save the modified Dog class in a script in IDLE and run it. Then open the interactive window and type the following to see instance methods in action:

```
>>> miles = Dog("Miles", 4)

>>> miles.description()
'Miles is 4 years old'

>>> miles.speak("Woof Woof")
'Miles says Woof Woof'

>>> miles.speak("Bow Wow")
'Miles says Bow Wow'
```

The `.description()` method defined in the above `Dog` class returns a string containing information about the `Dog` instance `miles`. When writing your own classes, it is a good idea to have a method that returns a string containing useful information about an instance of the class. However, `.description()` isn't the most Pythonic way of doing this.

When you create a `list` object, you can use the `print()` function to display a string that looks like the list:

```
>>> names = ["Fletcher", "David", "Dan"]
>>> print(names)
['Fletcher', 'David', 'Dan']
```

Let's see what happens when you `print()` the `miles` object:

```
>>> print(miles)
<__main__.Dog object at 0x00aeef70>
```

When you `print(miles)`, you get a cryptic looking message telling you that `miles` is a `Dog` object at `0x00aeef70`. The number `0x00aeef70` is the address of this `Dog` object in your computer's memory, and the number you see on your computer will be different.

The message displayed by `print(miles)` isn't very helpful. You can change what gets printed by defining a special instance method called `__str__()`.

Let's change `.description()` to `__str__()` in the `Dog` class:

```
class Dog:  
    # Leave other parts of Dog class as-is  
  
    # Replace .description() with __str__()  
    def __str__(self):  
        return f"{self.name} is {self.age} years old"
```

Now when you `print(miles)` you get much friendlier output:

```
>>> miles = Dog("Miles", 4)  
>>> print(miles)  
'Miles is 4 years old'
```

Note

Methods like `__str__()` are commonly called **dunder methods** because they begin and end with double underscores. There are a number of dunder methods available that allow your classes to work well with other Python language features.

Dunder methods are powerful and are an important part of mastering OOP in Python, but we won't go into detail here. For more information, you are encouraged to checkout [Operator and Function Overloading in Custom Python Classes](#) as well as Chapter 4 of [Python Tricks: The Book](#).

You should now have a pretty good idea of how to create a class that stores some data and provides some methods to interact with that data and define behaviors for an object.

In the next section, you'll see how to take your knowledge one step further and create classes from other classes. But first, check your

understanding with the following review exercises.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Modify the `Dog` class to include a third instance attribute called `coat_color` that stores the color of the dog's coat as a string. Store your new class in a script and test it out by adding the following code at the bottom of the script:

```
philo = Dog("Philo", 5, "brown")
print(f"{philo.name}'s coat is {philo.coat_color}.")
```

The output of your script should be:

Philo's coat is brown.

2. Create a `Car` class with two instance attributes: `.color`, which stores the name of the car's color as a string, and `.mileage`, which stores the number of miles on the car as an integer. Then instantiate two `Car` objects—a blue car with 20,000 miles, and a red car with 30,000 miles, and print out their colors and mileage. Your output should look like the following:

The blue car has 20,000 miles.

The red car has 30,000 miles.

3. Modify the `Car` class with an instance method called `.drive()` that takes a number as an argument and adds that number to the `.mileage` attribute. Test that your solution works by instantiating a car with 0 miles, then call `.drive(100)` and print the `.mileage` attribute to check that it is set to 100.

[Leave feedback on this section »](#)

10.3 Inherit From Other Classes

Inheritance is the process by which one class takes on the attributes and methods of another. Newly formed classes are called **child**

classes, and the classes that child classes are derived from are called **parent classes**.

Child classes can override and extend the attributes and methods of parent classes. In other words, child classes inherit all of the parent's attributes and methods but can also specify different attributes and methods that are unique to themselves, or even redefine methods from their parent class.

The concept of object inheritance can be thought of sort of like genetic inheritance, even though the analogy isn't perfect.

For example, you may have inherited your hair color from your mother. It's an attribute you were born with. You may decide that you want to color your hair purple. Assuming your mother doesn't have purple hair, you have just **overridden** the hair color attribute you inherited from your mom.

You also inherit, in a sense, your language from your parents. If your parents speak English, then you will also speak English. One day, you may decide to learn a second language, like German. In this case you are **extending** attributes, because you have added an attribute that your parents do not have.

The object Class

The most basic type of class is an `object`, which generally all other classes inherit from as their parent. When you define a new class, Python 3 implicitly uses `object` as the parent class, so the following two definitions are equivalent:

```
class Dog(object):  
    pass
```

In Python 3, this is the same as:

```
class Dog:
```

```
pass
```

The inheritance from `object` is stated explicitly in the first definition by putting `object` in between parentheses after the `Dog` class name. This is the same pattern used to create child classes from your own custom classes.

Note

In Python 2 there's a distinction between [new-style and old-style classes](#). We won't cover this distinction, because it doesn't apply to Python 3

Just know that in Python 3, there is an `object` class that all classes inherit from, even though you don't have to explicitly state that in your code.

Let's see how and why you might create child classes from a parent class.

Dog Park Example

Pretend for a moment that you are at a dog park. There are many dogs of different breeds at the park, all engaging in various dog behaviors.

Suppose now that you want to model the dog park with Python classes. The `Dog` class you wrote in the previous section can distinguish dogs by name and age, but not by breed.

You could modify the `Dog` class by adding a `.breed` attribute:

```
class Dog:  
    species = "Canis familiaris"  
  
    def __init__(self, name, age, breed):  
        self.name = name  
        self.age = age  
        self.breed = breed
```

The instance methods defined earlier are omitted here because they aren't important for this discussion.

Now, to model the dog park, you could instantiate a bunch of different dogs:

```
>>> miles = Dog("Miles", 4, "Jack Russell Terrier")
>>> buddy = Dog("Buddy", 9, "Dachshund")
>>> jack = Dog("Jack", 3, "Bulldog")
>>> jim = Dog("Jim", 5, "Bulldog")
```

Each breed of dog has slightly different behaviors. For example, bulldogs have a low bark that sounds like “woof” but dachshunds have a higher pitched bark that sounds more like “yap”.

Using just the `Dog` class, you must supply a string for the `sound` argument of the `.speak()` method every time you call it on a `Dog` instance:

```
>>> buddy.speak("Yap")
'Buddy says Yap'

>>> jim.speak("Woof")
'Jim says Woof'

>>> jack.speak("Woof")
'Jack says Woof'
```

Passing a string to every call to `.speak()` is repetitive and inconvenient. What's worse, the string representing the sound each `Dog` instance makes depends on the `.breed` attribute, but there is nothing stopping you, or someone using the `Dog` class you have created, from passing any string they wish.

You can simplify the experience of working with the `Dog` class by creating a child class for each breed of dog. This allows you to extend the functionality each child class inherits, including specifying a default argument for `.speak()`.

Parent Classes vs Child Classes

Let's create a child class for each of the three breeds mentioned above: Jack Russell Terrier, Dachshund, and Bulldog.

For reference, here is the full definition of the `Dog` class:

```
class Dog:  
    species = "Canis familiaris"  
  
    def __init__(self, name, age):  
        self.name = name  
        self.age = age  
  
    def __str__(self):  
        return f"{self.name} is {self.age} years old"  
  
    def speak(self, sound):  
        return f"{self.name} says {sound}"
```

Remember, to create a child class, you create new class with its own name and then put the name of the parent class in parentheses. The following creates three new child classes of the `Dog` class:

```
class JackRussellTerrier(Dog):  
    pass  
  
class Dachshund(Dog):  
    pass  
  
class Bulldog(Dog):  
    pass
```

With the child classes defined, you can now instantiate some dogs of specific breeds:

```
>>> miles = JackRussellTerrier("Miles", 4)  
>>> buddy = Dachshund("Buddy", 9)
```

```
>>> jack = Bulldog("Jack", 3)
>>> jim = Bulldog("Jim", 5)
```

Instances of child classes inherit all of the attributes and methods of the parent class:

```
>>> miles.species
'Canis familiaris'

>>> buddy.name
'Buddy'

>>> print(jack)
Jack is 3 years old

>>> jim.speak("Woof")
'Jim says Woof'
```

To determine which class a given object belongs to, you can use the built-in `type()` function:

```
>>> type(miles)
<class '__main__.JackRussellTerrier'>
```

What if you wanted to determine if `miles` is also an instance of the `Dog` class? You can do this with the built-in `isinstance()` function:

```
>>> isinstance(miles, Dog)
True
```

Notice that `isinstance()` takes two arguments, an object and a class. In the example above, `isinstance()` checks if `miles` is an instance of the `Dog` class and returns `True`.

All of the `miles`, `buddy`, `jack` and `jim` objects are instances of the `Dog` class, but `miles` is not an instance of the `Bulldog` class, and `jack` is not an instance of the `Dachshund` class:

```
>>> isinstance(miles, Bulldog)
False

>>> isinstance(jack, Dachshund)
False
```

More generally, all objects created from a child class are instances of the parent class, although they may not be instances of other child classes.

Now that you've got some child classes created for some different breeds of dogs, let's give each breed its own sound.

Extending the Functionality of a Parent Class

At this point, we have four classes floating around: a parent class—`Dog`—and three child classes—`JackRussellTerrier`, `Dachshund` and `Bulldog`. All three child classes inherit every attribute and method from the parent class, including the `.speak()` method.

Since different breeds of dogs have slightly different barks, we want to provide a default value for the `sound` argument of their respective `.speak()` methods. To do this, we need to override the `.speak()` method in the class definition for each breed. To override a method defined on the parent class, you define a method with the same name on the child class.

Let's see what this looks like for the `JackRussellTerrier` class:

```
class JackRussellTerrier(Dog):
    def speak(self, sound="Arf"):
        return f"{self.name} says {sound}"
```

The `.speak()` method is now defined on the `JackRussellTerrier` class with the default argument for `sound` set to "Arf". Now you can call `.speak()` and a `JackRussellTerrier` instance without passing an argument to `sound`:

```
>>> miles = JackRussellTerrier("Miles", 4)
>>> miles.speak()
'Miles says Arf'
```

Sometimes dogs make different barks, so if Miles gets angry and growls, you can still call `.speak()` with a different sound:

```
>>> miles.speak("Grrr")
'Miles says Grrr'
```

One advantage of class inheritance is that changes to the parent class will automatically propagate to their child classes. This occurs as long as the attribute or method being changed isn't overridden in the child class.

For example, let's say you decide to change the string returned by `.speak()` in the `Dog` class:

```
class Dog:
    # Other attributes and methods omitted...

    def speak(self, sound):
        return f"{self.name} barks: {sound}"
```

Now, when you create a new `Bulldog` instance named `jim`, the result of `jim.speak("Woof")` will be 'Jim barks: Woof' instead of 'Jim says Woof':

```
>>> jim = Bulldog("Jim", 5)
>>> jim.speak("Woof")
'Jim barks: Woof'
```

However, calling `.speak()` on a `JackRussellTerrier` instance won't show the new style of output:

```
>>> miles = JackRussellTerrier("Miles", 4)
>>> miles.speak()
'Miles says Arf'
```

Sometimes it makes sense to completely override a method from a par-

ent class. But in this instance, we don't want the `JackRussellTerrier` class to lose any changes that might be made to the formatting of the output string of `Dog.speak()`.

To do this, you still need to define a `.speak()` method on the `JackRussellTerrier` class. But instead of explicitly defining the output string, you need to call the `Dog` class's `.speak()` method *inside* of the child class's `.speak()` method and make sure to pass to it the whatever is passed to `sound` argument of `JackRussellTerrier.speak()`.

You can access the parent class from inside a method of a child class by using the `super()` function. Here's how you could re-write the `JackRussellTerrier.speak()` method using `super()`:

```
class JackRussellTerrier(Dog):
    def speak(self, sound="Arf"):
        return super().speak(sound)
```

When you call `super().speak(sound)` inside of `JackRussellTerrier`, Python searches the parent class, `Dog`, for a `.speak()` method and calls it with the variable `sound`. Now, when you call `miles.speak()`, you will see output reflecting the new formatting in the `Dog` class:

```
>>> miles = JackRussellTerrier("Miles", 4)
>>> miles.speak()
'Miles barks: Arf'
```

Important

In the above examples, the **class hierarchy** is very simple: the JackRussellTerrier class has a single parent class—Dog.

In many real world examples, the class hierarchy can get quite complicated with one class inheriting from a parent class, which inherits from another parent class, which inherits from another parent class, and so on.

The `super()` function does much more than just search the parent class for a method or an attribute. It traverses the entire class hierarchy for a matching method or attribute. If you aren't careful, `super()` can have surprising results.

In this section, you learned how to make new classes from existing classes utilizing an OOP concept called **inheritance**. You saw how to check if an object is an instance of a class or parent class using the `isinstance()` function. Finally, you learned how to extend the functionality of a parent class by using `super()`.

In the next section you will bring together everything you have learned by using classes to model a farm. Before you tackle the assignment, check your understanding with the review exercises below.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Create a GoldenRetriever class that inherits from the Dog class. Give the `sound` argument of the `GoldenRetriever.speak()` method a default value of "Bark". Use the following code for your parent Dog class:

```
class Dog:  
    species = "Canis familiaris"  
  
    def __init__(self, name, age):
```

```
self.name = name
self.age = age

def __str__(self):
    return f"{self.name} is {self.age} years old"

def speak(self, sound):
    return f"{self.name} says {sound}"
```

2. Write a `Rectangle` class that must be instantiated with two attributes: `length` and `width`. Add a `.area()` method to the class that returns the area (`length * width`) of the rectangle. Then write a `Square` class that inherits from the `Rectangle` class and that is instantiated with a single attribute called `side_length`. Test your `Square` class by instantiating a `Square` with a `side_length` of 4. Calling the `.area()` method should return 16.

[Leave feedback on this section »](#)

10.4 Challenge: Model a Farm

In this assignment, you'll create a simplified model of a farm. As you work through this assignment, keep in mind that there are a number of correct answers.

The focus of this assignment is less about the Python class syntax and more about software design in general, which is highly subjective. This assignment is intentionally left open-ended to encourage you to think about how you would organize your code into classes.

Before you write any code, grab a pen and paper and sketch out a model of your farm, identifying classes, attributes, and methods. Think about inheritance. How can you prevent code duplication? Take the time to work through as many iterations as you feel are necessary.

The actual requirements are open to interpretation, but try to adhere to these guidelines:

1. You should have at least four classes: the parent `Animal` class, and then at least three child animal classes that inherit from `Animal`.
2. Each class should have a few attributes and at least one method that models some behavior appropriate for a specific animal or all animals—such as walking, running, eating, sleeping, and so on.
3. Keep it simple. Utilize inheritance. Make sure you output details about the animals and their behaviors.

You can find the solutions to this code challenge and many other bonus resources online at realpython.com/python-basics/resources.

[Leave feedback on this section »](#)

10.5 Summary and Additional Resources

In this chapter you learned about object-oriented programming (OOP) in Python, which is a programming paradigm that is not specific to Python. Most of the modern programming languages—such as Java, C#, and C++—follow OOP principles.

You saw how to define a class, which is a sort of “blueprint” for an object, and how to instantiate an object from a class. You also learned about attributes, which correspond to properties of an object, and methods, which correspond to behaviors and actions of an object.

Finally, you learned how inheritance works by creating child classes from a parent class. You saw how to reference a method on a parent class using `super()`, and how to check if an object inherits from some class using `isinstance()`.

OOP is a big and sometimes difficult topic. Some programmers consider OOP a foundational part of modern programming, but this viewpoint isn’t without its [criticisms](#).

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-10

Additional Resources

You've seen the basics of OOP, but there is so much more to learn! Continue your journey with the following resources:

1. [Official Python documentation](#)
2. [OOP Articles on Real Python](#)
3. [Recommended resources on realpython.com](#)

[Leave feedback on this section »](#)

Chapter 11

Modules and Packages

As you gain experience writing code, you will eventually work on projects that are so large that keeping all of the code in a single file becomes cumbersome.

Instead of writing a single file, you can put related code into separate files called **modules**. Individual modules can be put together like building blocks to create a larger application.

There are four main advantages to breaking a program into modules:

1. **Simplicity:** Modules are focused on a single problem.
2. **Maintainability:** Small files are better than large files.
3. **Reusability:** Modules reduce duplicate code.
4. **Scoping:** Modules have their own **namespaces**.

In this chapter, you will learn how to:

- Create your own modules
- Use modules in another file via the **import** statement
- Organize several modules into a **package**

Let's get started!

[Leave feedback on this section »](#)

11.1 Working With Modules

A **module** is a file containing Python code that can be re-used in other Python code files.

Technically, every Python script file that you have created while reading this book is a module, but you haven't seen how to use code from one module inside of another.

In this section, you'll explore modules in more detail. You'll learn how to create them with IDLE, how to import one module into another, and understand how modules create namespaces.

Creating Modules

Open IDLE and start a new script window by selecting **File** \gg **New File** or by pressing **Ctrl** + **N**. In the script window, define a function `add()` that returns the sum of its two parameters:

```
# adder.py

def add(x, y):
    return x + y
```

Select **File** \gg **Save** or press **Ctrl** + **S** to save the file as `adder.py` in a new directory called `myproject/` somewhere on your computer. `adder.py` is a Python module! It's not a complete program, but not all modules need to be.

Now open another new script window by pressing **Ctrl** + **N** and type the following code:

```
# main.py

value = add(2, 2)
```

```
print(value)
```

Save the file as `main.py` in the same `myproject/` folder you just created. Then press **F5** to run the module.

When the module runs you'll see a `NameError` displayed in IDLE's interactive window:

```
Traceback (most recent call last):
  File "//Documents/myproject/main.py", line 1, in <module>
    value = add(2, 2)
NameError: name 'add' is not defined
```

It makes sense that a `NameError` occurs because `add()` is defined in `adder.py` and not in `main.py`. In order to use `add()` in `main.py`, you must first import the `adder` module.

Importing One Module Into Another

In the script window for `main.py`, add the following line to the top of the file:

```
# main.py

import adder # <-- Add this line

# Leave the code below unchanged
value = add(2, 2)
print(value)
```

When you `import` one module into another, the contents of the imported module become available in the other. The module with the `import` statement is called the **calling module**. In this example, `adder.py` is the imported module and `main.py` is the calling module.

Press **Ctrl**+**S** to save `main.py` and press **F5** to run the module. The `NameError` exception is still raised. That's because `add()` can only be accessed from the `adder` namespace.

A **namespace** is a collection of names, such as variable names, function names, and class names. Every Python module has its own namespace.

Variables, functions, and classes in a module can be accessed from within the same module by just typing their name. That's how you've been doing it throughout this book so far. However, this doesn't work for imported modules.

To access a name in an imported module from the calling module, type the imported module's name followed by a dot (.) and the name you want to use:

```
<module>. <name>
```

For instance, to use the `add()` function in the `adder` module, you need to type `adder.add()`.

Important

The name used to import a module is the same as the module's file name.

For this reason, module file names must be valid Python identifiers. That means they may only contain upper and lower case letters, numbers, and underscores (_), and they may not start with a digit.

Now update the code in `main.py` as follows:

```
# main.py

import adder

value = adder.add(2, 2) # <-- Change this line
print(value)
```

Save the file and run the module. The value 4 is printed in the interactive window.

When you type `import <module>` at the beginning of a file, the module's entire namespace is imported. Any new variables or functions added to `adder.py` will be accessible in `main.py` without having to import anything new.

Open the script window for `adder.py` and add the following function below `add()`:

```
# adder.py

# Leave this code unchanged
def add(x, y):
    return x + y

def double(x): # <-- Add this function
    return x + x
```

Save the file. Then open the script window for `main.py` and add the following code:

```
# main.py

import adder

value = adder.add(2, 2)
double_value = adder.double(value) # <-- Add this line
print(double_value) # <-- Change this line
```

Now save and run `main.py`. When the module runs, the value 8 is displayed in the interactive window. Since `double()` already exists in the `adder` namespace, no `NameError` is raised.

Import Statement Variations

The `import` statement is flexible. There are two variations that you should know about:

1. `import <module> as <other_name>`

2. `from <module> import <name>`

Let's look at each of these variations in detail.

```
import <module> as <other_name>
```

You can change the name of an import using the `as` keyword:

```
import <module> as <other_name>
```

When you import a module this way, the module's namespace is accessed through `<other_name>` instead of `<module>`.

For example, change the `import` statement in `main.py` to the following:

```
import adder as a # <-- Change this line

# Leave the code below unchanged
value = adder.add(2, 2)
double_value = adder.double(value)
print(double_value)
```

Save the file and press `F5`. A `NameError` is raised:

```
Traceback (most recent call last):
  File "//Mac/Home/Documents/myproject/main.py", line 3, in <module>
    value = adder.add(2, 2)
NameError: name 'adder' is not defined
```

The `adder` name is no longer recognized because the module has been imported with the name `a` instead of `adder`.

To make `main.py` work, you need to replace `adder.add()` and `adder.double()` with `a.add()` and `a.double()`:

```
import adder as a

value = a.add(2, 2) # <-- Change this line
```

```
double_value = a.double(value) # <-- Change this line, too
print(double_value)
```

Now save the file and run the module. No `NameError` is raised and the value 8 is displayed in the interactive window.

```
from <module> import <name>
```

Instead of importing the entire namespace, you can import only a specific name from a module. To do this, replace the `import` statement with the following:

```
from <module> import <name>
```

For example, in `main.py`, change the `import` statement to the following:

```
from adder import add # <-- Change this line

value = adder.add(2, 2)
double_value = adder.double(2, 2)
print(double_value)
```

Save the file and press `F5`. A `NameError` exception is raised:

```
Traceback (most recent call last):
  File "//Documents/myproject/main.py", line 3, in <module>
    value = adder.add(2, 2)
NameError: name 'adder' is not defined
```

The above traceback tells you that the name `adder` is undefined. Only the name `add` is imported from `adder.py` and is placed in the `main.py` module's local namespace. That means you can use `add()` without having to type `adder.add()`.

Replace `adder.add()` and `adder.double()` in `main.py` with `add()` and `double()`:

```
from adder import add

value = add(2, 2) # <-- Change this line
double_value = double(value) # <-- Change this line, too
print(double_value)
```

Now save the file and run the module. What do you think happens?

Another `NameError` is raised:

```
Traceback (most recent call last):
  File "//Documents/myproject/main.py", line 4, in <module>
    double_value = double(value)
NameError: name 'double' is not defined
```

This time, the `NameError` tells you that the name `double` isn't defined, which proves that only the `add` name was imported from the `adder` module.

You can import the `double` name by adding it to the import statement in `main.py`:

```
from adder import add, double # <-- Change this line

# Leave the code below unchanged
value = add(2, 2)
double_value = double(value)
print(double_value)
```

Save and run the module. Now the module runs without producing a `NameError`. The value `8` is displayed in the interactive window.

Summary of Import Statements

The following table summarizes what you've learned about importing modules:

Import Statement	Result
<code>import <module></code>	Import all of <code><module></code> 's namespace into the name <code><module></code> . Import module names can be accessed from the calling module with <code><module>.name</code> .
<code>import <module> as <other_name></code>	Import all of <code><modules></code> 's namespace into the name <code><other_name></code> . Import module names can be accessed from the calling module with <code><other_name>.name</code> .
<code>from <module> import <name1>, <name2>, ...</code>	Import only the names <code><name1></code> , <code><name2></code> , etc, from <code><module></code> . The names are added to the calling modules's local namespace and can be accessed directly.

Separate namespaces are one of the great advantages of dividing code into individual modules, so let's take some time to explore why namespaces matter and why you should care about them.

Why Use Namespaces?

Suppose every person on the entire planet is given an ID number. In order to distinguish one person from the next, each ID number needs to be unique. We'll need a whole bunch of ID numbers to make that work!

The world is divided into countries, and we can group people by their country of birth. If we assign each country a unique code, we can include that code in a person's ID number. For example, a person from the United States might have an ID of `US-357`, and a person from Great Britain might have an ID of `GB-246`.

Two people from different countries can now have the same ID number. We can distinguish them because their ID's begin with different country codes. Every person from the same country must have a unique ID number, but we no longer need globally unique ID numbers.

The country codes in this scenario are an examples of **namespaces**, and illustrate three of the main reasons namespaces are used:

1. They **group** names into logical containers
2. They **prevent clashes** between duplicate names
3. They **provide context** to names

Namespace in code provide the same advantages.

You have seen three different ways to import a module into another one. Keeping in mind the advantages namespaces give you can help you determine which kind of `import` statement makes the most sense.

In general, `import <module>` is the preferred way to import a module because it keeps the imported module's namespace completely separate from the calling module's namespace. Moreover, every name from the imported module is accessed from the calling module with the `<module>. <name>` format, which immediately tells you which module the name originates in.

There are two reasons you might use the `import <module> as <other_name>` format:

1. The module name is long and you wish to import an abbreviated version of it
2. The module name clashes with an existing name in the calling module

`import <module> as <other_name>` still keeps the imported module's namespace separate from the calling module's namespace. The tradeoff is that the name you give the module might not be as easily recognizable as the original module name.

Importing specific names from a module is generally the least preferred way to import code from a module. The imported names are added directly to the calling module's namespace, completely removing them from context of the calling module.

Sometimes, modules contain a single function or class that has the same name as the module. For example, there is a module in the Python standard library called `datetime` that contains a class called `datetime`.

Suppose you add the following `import` statement to your code:

```
import datetime
```

This imports the `datetime` module into your code's namespace, so in order to use the `datetime` class contained in the `datetime` module, you need to type the following:

```
datetime.datetime(2020, 2, 2)
```

Don't worry about how the `datetime` class works right now. The important part of this example is that having to constantly type `datetime.datetime` anytime you want to use the `datetime` class is redundant and tiring.

This is a great example of when it's appropriate to use one of the variations of the `import` statement. To keep the context of the `datetime` package, it is common for Python programmers to import the package and rename it as `dt`:

```
import datetime as dt
```

Now, to use the `datetime` class, you only need to type `dt.datetime`:

```
dt.datetime(2020, 2, 2)
```

It is also common for Python programmers to import the `datetime` class directly into the calling module's namespace:

```
from datetime import datetime
```

This is fine because the context isn't really lost. The class and the module share the same name, after all.

When imported directly, you no longer have to use dotted module

names to access the `datetime` class:

```
datetime(2020, 2, 2)
```

The various import statements allow you to reduce typing and unnecessarily long dotted module names. That said, abusing the various `import` statements can lead to a loss of context, resulting in code that is more difficult to understand.

Always use good judgment when importing modules so that the most context possible is preserved.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Create a module called `greeter.py` that contains a single function `greet()`. This function should accept a single string parameter `name` and print the text `Hello {name}!` to the interactive window with `{name}` replaced with the function argument.
2. Create a module called `main.py` that imports the `greet()` function from `greeter.py` and calls the function with the argument "Real Python".

[Leave feedback on this section »](#)

11.2 Working With Packages

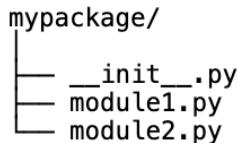
Modules allow you to divide a program into individual files that can be reused as needed. Related code can be organized into a single module and kept separate from other code.

Package take this organizational structure one step further by allowing you to group related modules under a single namespace.

In this section, you'll learn how to create your own Python package and import code from that package into another module.

Creating Packages

A **package** is a folder that contains one or more Python modules. It must also contain a special module called `__init__.py`. Here is an example of a package so that you can see this structure:



The `__init__.py` module doesn't need to contain any code! It only needs to exist so that Python recognizes the `mypackage/` folder as a Python package.

Using your computers file explorer, or whatever tool you are comfortable with, create a new folder somewhere on your computer called `packages_example/`. Inside of that folder, create another folder called `mypackage/`.

The `packages_example/` folder is called the **project folder**, or **project root folder**, because it contains all of the files or folders in the `packages_examples` project. The `mypackage/` folder will eventually become a Python package. It isn't one right now because it doesn't contain any modules.

Open IDLE and create a new script window by pressing **Ctrl + N**. At the top of the file, add the following comment:

```
# main.py
```

Now press **Ctrl + S** and save the file as `main.py` in the `packages_example/` folder you created earlier.

Open another script window by pressing **Ctrl + N**. Insert the following at the top of the file:

```
# __init__.py
```

Then save the file as `__init__.py` in the `mypackage/` subfolder of your `packages_example` folder.

Finally, create two more script windows. Save these files as `module1.py` and `module2.py`, respectively, in your `mypackages/` folder, and insert comments at the top of each file containing the file name.

When you are done, you should have five IDLE windows open: the interactive window and four script windows. You can arrange your screen to look something like this:

11.2. Working With Packages

The screenshot shows a code editor with five windows open:

- __init__.py**: Contains the code `# __init__.py`.
- module1.py**: Contains the code `def greet(name): print(f"Hello, {name}!")`.
- module2.py**: Contains the code `def depart(name): print(f"Goodbye, {name}!")`.
- main.py**: Contains the code `# main.py`.
- Python 3.8.1 Shell**: Displays the Python environment information and a prompt `>>>`.

Now that we've created the package structure, let's add some code. In the `module1.py` file, add the following function:

```
# module1.py

def greet(name):
    print(f"Hello, {name}!")
```

In the `module2.py` file add the following:

```
def depart(name):
    print(f"Goodbye, {name}!")
```

Make sure you save both of the `module1.py` and `module2.py` files! You're now ready to import and use these modules in the `main.py` module.

Importing Modules From Packages

In your `main.py` file, add the following code:

```
# main.py

import mypackage

mypackage.module1.greet("Pythonista")
mypackage.module2.depart("Pythonista")
```

Save `main.py` and press `F5` to run the module. In the interactive window, an `AttributeError` is raised:

```
Traceback (most recent call last):
  File "\MacHome\Documents\packages_examplemain.py", line 5, in <module>
    mypackage.module1.greet("Pythonista")
AttributeError: module 'mypackage' has no attribute 'module1'
```

When you import the `mypackage` module, the `module1` and `module2` namespaces are not imported automatically. In order to use them, you need to import them as well.

Change the import statement at the top of the `main.py`:

```
# main.py

import mypackage.module1 # <-- Change this line

# Leave the below code unchanged
mypackage.module1.greet("Pythonista")
```

```
mypackage.module2.depart("Pythonista")
```

Now save and run the `main.py` module. You should see the following output in the interactive window:

```
Hello, Pythonista!  
Traceback (most recent call last):  
  File "\MacHomeDocuments\packages_exampemain.py", line 6, in <module>  
    mypackage.module2.depart("Pythonista")  
AttributeError: module 'mypackage' has no attribute 'module2'
```

You can tell that `mypackage.module1.greet()` was called because `Hello, Pythonista!` is displayed in the interactive window.

However, `mypackage.module2.depart()` was not called. That line raised an attribute error because the only module imported from `mypackage` so far is `module1`.

To import `module2`, add the following `import` statement to the top of your `main.py` file:

```
# main.py  
  
import mypackage.module1  
import mypackage.module2 # <-- Add this line  
  
# Leave the below code unchanged  
mypackage.module1.greet("Pythonista")  
mypackage.module2.depart("Pythonista")
```

Now when you save and run `main.py`, both `greet()` and `depart()` get called:

```
Hello, Pythonista!  
Goodbye, Pythonista!
```

In general, modules are imported from packages using **dotted module names** with the following format:

```
import <package_name>.<module_name>
```

First type the name of the package followed by a dot (.) and the name of the module you want to import.

Important

Just like module file names, package folder names must be valid Python identifiers. They may only contain upper and lower case letters, numbers, and underscores (_), and they may not start with a digit.

As with modules, there are several variations on the `import` statement that you can use when importing packages.

Import Statement Variations For Packages

There are three variations of the `import` statement that you learned for importing names from modules. These three variations translate to the following four variations for importing modules from packages:

1. `import <package>`
2. `import <package> as <other_name>`
3. `from <package> import <module>`
4. `from <package> import <module> as <other_name>`

These variations work much the same was as the counterparts for modules.

For instance, rather than importing `mypackage.module1` and `mypackage.module2`, you can import both modules from the package on the same line.

Change your `main.py` file to the following:

```
# main.py  
  
from mypackage import module1, module2
```

```
module1.greet("Pythonista")
module2.depart("Pythonista")
```

When you save and run the module, the same output as before is displayed in the interactive window.

You can change the name of an imported module using the `as` keyword:

```
# main.py

from mypackage import module1 as m1, module2 as m2

m1.greet("Pythonista")
m2.depart("Pythonista")
```

You can also import individual names from a package module. For instance, you can rewrite your `main.py` to the following without changing what gets printed when you save and run the module:

```
# main.py

from mypackage.module1 import greet
from mypackage.module2 import depart

greet("Pythonista")
depart("Pythonista")
```

With so many ways to import packages, it's natural to wonder which way is best.

Guidelines For Importing Packages

The same guidelines for importing names from modules apply to importing modules from packages. You should prefer that imports be as explicit as possible, so that the modules and names imported from the package into the calling module have the appropriate context.

In general, the following format is the most explicit:

```
import <package>.<module>
```

Then, to access names from module, you need to type something like the following

```
<package>.<module>.<name>
```

When you encounter names that are used from the imported module, there is no question where those names come from. But sometimes package and module names are long, and you find yourself typing `<package>.<module>` over and over again in your code.

The following format allows you to skip the package name and import just the module name into the calling module's namespace:

```
from <package> import <module>
```

Now you can just type `<module>.<name>` to access some name from the module. While this no longer tells you from which package the name comes from, it does keep the context of the module apparent.

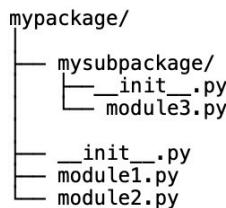
Finally, the following format is generally ambiguous and should only be used when there is no risk of importing a name from a module that clashes with a name in the calling module:

```
from <package>.<module> import <name>
```

Now that you've seen how to import modules from packages, let's take a quick look at how to nest packages inside of other packages.

Importing Modules From Subpackages

A package is just a folder containing one or more Python modules, one of which must be named `__init__.py`, so it's entirely possible to have the following package structure:



A package nested inside of another package is called a **subpackage**. For example, the `mysubpackage` folder is a subpackage of `mypackage` because it contains an `__init__.py` module, as well as a second module called `module3.py`.

Using your computer's file explorer, or some other tool, create the `mysubpackage/` folder on your computer. Make sure you place the folder inside of the `mypackage/` folder you created earlier.

In IDLE, open two new script windows. Create the files `__init__.py` and `module3.py` and save both modules to the `mysubpackage/` folder.

In your `module3.py` file, add the following code:

```
# module3.py

people = ["John", "Paul", "George", "Ringo"]
```

Now open the `main.py` file in your root `packages_examples/` project folder. Remove any existing code and replace it with the following:

```
# main.py

from mypackage.module1 import greet
from mypackage.mysubpackage.module3 import people

for person in people:
    greet(person)
```

The `people` list from the `module3` module inside of `mysubpackage` is im-

ported via the dotted module name `mypackage.mysubpackage.module3`.

Now save and run `main.py`. The following output is displayed in the interactive window:

```
Hello, John!  
Hello, Paul!  
Hello, George!  
Hello, Ringo!
```

Subpackages are great for organizing code inside of very large packages. They help keep the folder structure of a package clean and organized.

However, deeply nested subpackages introduce long dotted module names. You can image how much typing it would take to import a module from a subpackage of a subpackage of a subpackage of a package.

It's good practice to try and keep your subpackages at most one or two levels deep.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. In a new project folder called `package_exercises/`, create a package called `helpers` with three modules: `__init__.py`, `string.py`, and `math.py`.

In the `string.py` module, add a function called `shout()` that takes a single string parameter and returns a new string with all of the letters in uppercase.

In the `math.py` module, as a function called `area()` that takes two parameters called `length` and `width` and returns their product `length * width`.

2. In the root project folder, create a module called `main.py` that imports the `shout()` and `area()` functions. Use the `shout()` and `area()` functions to print the following output:

```
THE AREA OF A 5-BY-8 RECTANGLE IS 40
```

[Leave feedback on this section »](#)

11.3 Summary and Additional Resources

In this chapter you learned how to create your own Python modules and packages, and how to import objects from one module into another.

You saw that dividing code into modules and packages is advantageous because:

- Small code files are **simpler** than large code files
- Small code files are **easier to maintain** than large code files
- Modules can be **reused** throughout a project
- Modules group related objects together into isolated **namespaces**

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-modules-and-packages

Additional Resources

To learn more about modules and packages, check out the following resources:

- Python Modules and Packages Course
- Absolute and Relative Imports
- Recommended resources on realpython.com

[Leave feedback on this section »](#)

Chapter 12

File Input and Output

So far, you've written programs that get their input from one of two places: from the program itself or from the user. Program output has been limited to displaying some text in IDLE's interactive window.

These input and output methods are not useful in several common scenarios:

- The input values are unknown while writing the program
- The program requires more data than a user can be expected to type in by themselves
- Output must be shared with other people after the program runs

This is where files come in.

In this chapter, you will learn how to:

- Work with file paths and file metadata
- How to read and write text files
- How to read and write **Comma-Separated Value (CSV)** files
- How to create, delete, copy, and move files and folders

Let's dive in!

[Leave feedback on this section »](#)

12.1 Files and the File System

You have likely been working with computer files for a long time. Even so, there are some things that programmers need to know about files that the general user does not.

In this section, you'll learn the concepts necessary to get started working with files in Python.

Note

If you are familiar with concepts like the file system and file paths, may wish to read the *Working With File Paths in Python* and *File Metadata* sections before skipping to the next section.

Let's start by exploring what a file is and how computers interact with them.

The Anatomy of a File

There are a multitude of types of files out there: text files, image files, audio files, and PDF files, just to name a few. Ultimately, though, a file is just a sequence of **bytes** called the **contents** of the file.

Each byte in a file can be thought of as an integer with a value between 0 and 255, including both endpoints. The bytes are the values that are stored on a physical storage device when a file is saved.

When you access a file on a computer, the contents of the file are read from the disk in the correct sequence of bytes. The important thing to know here is that there is nothing intrinsic to the file itself that dictates how to interpret the contents.

As a programmer, it's your job to properly interpret the contents when you open a file. This might sound difficult, but Python does a lot of the hard work for you.

For example, when you open a text file, Python can convert the numerical bytes of the file into text characters for you. You do not need to know the specifics of how this conversion happens. There are tools in the standard library for working with all sorts of file types, including images and audio files.

In order to access a file from a storage device, a whole host of things need to happen. You need to know on which device the file is stored, how to interact with that device, and where exactly on the device the file is located.

This monumental task is managed by a file system. Python interacts with the file system on your computer in order to read, write, and manipulate files.

The File System

The file system on a computer does two things:

1. It provides an abstract representation of the files stored on your computer and devices connected to it.
2. It interfaces with devices to control storage and retrieval file data.

Python interacts with the file system on your computer, so you can only do in Python whatever your file system allows.

Important

Different operating systems use different file systems. This is very important to keep in mind when writing code that will be run on different operating systems.

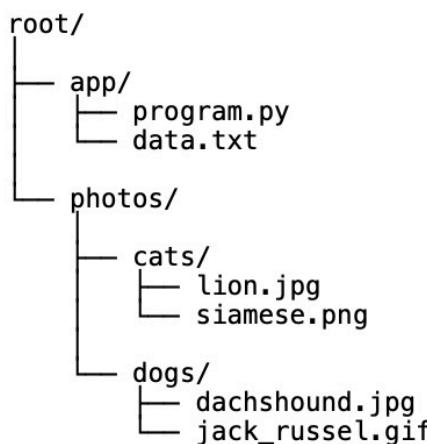
The file system itself manages communication between the computer and the physical storage device, so the only part of the file system you need to understand as a programmer is how it represents files.

The File System Hierarchy

File systems organize files in a hierarchy of **directories**, which are also known as **folders**. At the top of the hierarchy is a directory called the **root directory**. All other files and directories in the file system are contained in the root directory.

Each file in directory has a **file name** that must be unique from any other file in the same directory. Directories can also contain other directories, called **subdirectories** or **subfolders**.

The following **directory tree** visualizes the hierarchy of files and directories in an example file system:



In this file system, the root folder is called `root/`. It has two subdirectories: `app/` and `photos/`. The `app/` subdirectory contains a `program.py` file and a `data.txt` file. The `photos/` directory also has two subdirectories, `cats/` and `dogs/`, that both contains two image files.

File Paths

To locate a file in a file system, you can list the directories in order, starting with the root directory, followed by the name of the file. A

string with the file location represented in this manner is called a **file path**.

For example, the file path for the `jack_russel.gif` file in the above file system is `root/photos/dogs/jack_russel.gif`.

How you write file paths depends on your operating system. Here are three examples of file paths on Windows, macOS, and Linux:

1. **Windows:** `C:\Users\David\Documents\hello.txt`
2. **macOS:** `/Users/David/Documents/hello.txt`
3. **Ubuntu Linux:** `/home/David/Documents/hello.txt`

All three of these file paths locate a text file named `hello.txt` that is stored in the `Documents` subfolder of the user directory for a user named David. As you can see, there are some pretty big differences between file paths from one operating system to another.

On macOS and Ubuntu Linux, the operating system uses a **virtual file system** that organizes all files and directories for all devices on the system under a single root directory, usually represented by a forward slash symbol (`/`). Files and folders from external storage devices are usually located in a subdirectory called `media/`.

In Windows, there is no universal root directory. Each device has a separate file system with a unique root directory that is named with a **drive letter** followed by a colon (`:`) and a back slash symbol (`\`). Typically, the hard drive where the operating system is installed is assigned the letter `c`, so the root directory of the file system for that drive is `c:\`.

The other major difference between Windows, macOS, and Ubuntu file paths is that directories in a Windows file path are separated by back slashes (`\`), whereas directories in macOS and Ubuntu file paths are separated by forward slashes (`/`).

When you write programs that need to run on multiple operating systems, it is critical that you handle the differences in file paths appro-

priately. In versions of Python greater than 3.4, the standard library contains a module called `pathlib` helps take the pain out of handling file paths across operating systems.

Read on to learn how to use `pathlib` to work with file paths in Python.

[Leave feedback on this section »](#)

12.2 Working With File Paths in Python

To work with file paths in Python, use the standard libraries `pathlib` module. You'll need to import the module before you can do anything with it.

Open IDLE's interactive window and type the following to import the `pathlib` module:

```
>>> import pathlib
```

The `pathlib` module contains a class called `Path` that is used to represent a file path.

Creating Path Objects

There are several ways to create a new `Path` object:

1. From a string
2. With `Path.home()` and `Path.cwd()` class methods
3. With the `/` operator

The most straightforward way to create a `Path` object is from a string.

Creating Path Objects from Strings

For instance, the following creates a `Path` object representing the macOS file path `"/Users/David/Documents/hello.txt"`:

```
>>> path = pathlib.Path("/Users/David/Documents/hello.txt")
```

There's problem, though, with Windows paths. On Windows, directories are separated by back slashes \. Python interprets back slashes as the start of an **escape sequence** that represent a special character in the string, such as the newline character (\n).

Attempting to create a Path object with the Windows file path "C:\Users\David\Desktop\hello.txt" raises an exception:

```
>>> path = pathlib.Path("C:\Users\David\Desktop\hello.txt")
SyntaxError: (unicode error) 'unicodeescape' codec can't decode bytes
in position 2-3: truncated \UXXXXXXXXX escape
```

There are two ways to get around this problem:

You can use a forward slash (/) instead of a back slash (\) in your Windows file paths

```
>>> path = pathlib.Path("C:/Users/David/Desktop/hello.txt")
```

Python can interpret this just fine and will translate the path appropriately and automatically when interfacing with the Windows operating system.

You can also turn the string into a raw string by prefixing it with an r:

```
>>> path = pathlib.Path(r"C:\Users\David\Desktop\hello.txt")
```

This tells Python to ignore any escape sequences and just read the string as-is.

Path.home() and Path.cwd()

Besides creating a Path object from a string, the Path class has class methods that return Path objects of special directories. Two of the most useful class methods are Path.home() and Path.cwd().

Every operating system has a special directory for storing data for the

currently logged in user. This directory is called the user's **home directory**. The location of this directory depends on the operating system:

- **Windows:** C:\Users<username>
- **macOS:** /Users/<username>
- **Ubuntu Linux:** /home/<username>

The `Path.home()` class method creates a `Path` object representing the home directory regardless of which operating system the code runs on:

```
>>> home = pathlib.Path.home()
```

When you inspect the `home` variable on Windows, you will see something like this:

```
>>> home  
WindowsPath("C:\\\\Users\\\\David")
```

The `Path` object created is a subclass of `Path` called `WindowsPath`. On other operating systems, the `Path` object returned is a subclass called `PosixPath`.

For example, on macOS, inspecting `home` will display something like the following:

```
>>> home  
PosixPath("/Users/David")
```

For the rest of this section, `WindowsPath` objects will be shown in the example output. However, all of the examples will work with `PosixPath` objects.

Note

`WindowsPath` and `PosixPath` objects share the same methods and attributes. From a programming standpoint, there is no difference between the two types of `Path` objects.

The `Path.cwd()` class method returns a `Path` object representing the **current working directory**, or **CWD**. The current working directory is a dynamic reference to a directory that depends on where a process on the computer is currently working.

When you run IDLE, the current working directory is usually set to the `Documents` directory in the current user's home directory:

```
>>> pathlib.Path.cwd()
WindowsPath(r"C:\Users\David\Documents")
```

This is not always the case, though. Moreover, the current working directory may change during the lifetime of a program.

`Path.cwd()` is useful, but be careful when you use it. When you do, make sure you know that the current working directory refers the directory that you expect it to.

Using the / Operator

If you have an existing `Path` object, you can use the `/` operator to extend the path with subdirectories or file names.

For example, the following creates a `Path` object representing a file named `hello.txt` in the `Documents` subdirectory of the current user's home directory:

```
>>> home / "Desktop" / "hello.txt"
WindowsPath('C:/Users/David/Desktop/hello.txt')
```

The `/` operator must always have a `Path` object on the left hand side. The right hand side can have either string representing a single file or directory, or a string representing a path, or another `Path` object.

Absolute vs. Relative Paths

A path that begins with the root directory in a file system is called an **absolute file path**. Not all file paths are absolute. A file path that is not absolute is called a **relative file path**.

Here's an example of a `Path` object that references a relative path:

```
>>> # Relative Windows path
>>> path = pathlib.Path(r"Photos\image.jpg")

>>> # Relative macOS or Linux path
>>> path = pathlib.Path("Photos/image.jpg")
```

Notice that the path string does not start with `c:\` on Windows, or `/` on macOS and Linux.

You can check whether or not a file path is absolute using the `.is_absolute()` method:

```
>>> path.is_absolute()
False
```

Relative paths only make sense when considered within the context of some other directory. They are perhaps most commonly used to describe the path to a file relative to the current working directory, or the user's home directory.

You can extend a relative path to an absolute path using the forward slash (`/`) operator:

```
>>> home = pathlib.Path.home()
WindowsPath('C:/Users/David')
>>> home / pathlib.Path(r"Photos\image.png")
WindowsPath('C:/Users/David/Photos/image.png')
```

On the left of the forward slash (`/`), put an absolute path to the directory that contains the relative path. Then put the relative path on the right side of the forward slash.

Once you create a `Path` object, you can inspect the various components of the file path that it refers to.

Accessing File Path Components

All file paths contain a list of directories. The `.parents` attribute of a `Path` object returns an iterable containing the list of directories in the file path:

```
>>> path = pathlib.Path.home() / "hello.txt"
>>> path
WindowsPath("C:\\\\Users\\\\David")
>>> list(path.parents)
[WindowsPath("C:\\\\Users\\\\David"), WindowsPath("C:\\\\Users"),
WindowsPath("C:\\\\")]
```

Notice that the list of the directories are returned in reverse order from how they appear in the file path. That is, the last directory in the path is the first directory in the list of parent directories.

You can iterate over the parent directories in a `for` loop:

```
>>> for directory in path.parents:
...     print(directory)
...
C:\\Users\\David
C:\\Users
C:\\
```

The `.parent` attribute returns the name of the first parent directory in the file path as a string:

```
>>> path.parent
'C:\\Users\\David'
```

If the file path is absolute, you can access the root directory of the file path with the `.anchor` attribute:

```
>>> path.anchor  
'C:\'
```

Note that `.anchor` returns a string, and not another `Path` object.

For relative paths, `.anchor` return an empty string:

```
>>> path = pathlib.Path("hello.txt")  
>>> path.anchor  
''
```

The `.name` attribute returns the name of the file or directory that the path points to:

```
>>> home = pathlib.Path.home()      # C:\Users\David  
>>> home.name  
'David'  
>>> path = home / "hello.txt"  
>>> path.name  
'hello.txt'
```

The name of a file is broken down into two parts. The part to the left of the dot (.) is called the **stem**, and the part to the right of the dot (.) is called the **suffix** or **file extension**.

The `.stem` and `.suffix` attributes return strings containing each of these parts of the file name:

```
>>> path = pathlib.Path.home() / "hello.txt"  
>>> path.stem  
'hello'  
>>> path.suffix  
'.txt'
```

You might be wondering at this point how to actually do something with the `hello.txt` file. You'll learn how to read and write files in the next section. But before you open a file for reading, it might be a good idea to know whether or not that file exists.

Checking Whether Or Not a File Path Exists

You can create a `Path` object for a file path even if that path doesn't actually exist. Of course, file paths that don't represent actual files or directories aren't very useful, unless you plan on creating them at some point.

`Path` objects have an `.exists()` method that returns `True` or `False` depending on whether or not the file path exists on the machine executing the program.

For instance, if you don't have a `hello.txt` file in your home directory, then the `.exists()` method on the `Path` object representing that file path returns `False`:

```
>>> path = pathlib.Path.home() / "hello.txt"
>>> path.exists()
False
```

Using a text editor, or some other means, create a blank text file called `hello.txt` in your home directory. Then re-run the code from above, making sure `path.exists` returns `True`.

You can check whether or not a file path refers to a file or a directory. To check if the path is a file, use the `.is_file()` method:

```
>>> path.is_file()
True
```

Note that if the file path refers to a file, but doesn't exist, then `.is_file()` returns `False`.

Use the `.is_dir()` method to check if the file path refers to a directory

```
>>> # The path to "hello.txt" is not a directory
>>> path.is_dir()
False

>>> # The path to the home directory is a directory
```

```
>>> home.is_dir()  
True
```

Working with file paths is an essential part of any programming project that reads or writes data from a hard drive or other storage device. Understanding the differences between file paths on different operating systems and how to work with `pathlib.Path` objects so that your programs can work on any operating system is an important and useful skill.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Create a new `Path` object to a file called `my_file.txt` in a folder called `my_folder/` in your computer's home directory. Assign this `Path` object to the variable name `file_path`.
2. Check whether or not the path assigned to `file_path` exists.
3. Print the name of the path assigned to `file_path`. The output should be `my_file.txt`.
4. Print the name of the parent directory of the path assigned to `file_path`. The output should be `my_folder`.

[Leave feedback on this section »](#)

12.3 Common File System Operations

Now that you have a good grasp on the file system and working with file paths using the `pathlib` module, let's take a look at some common file operations and how you do them in Python.

In this section, you'll learn how to:

- Create directories and files
- Iterate over the contents of a directory

- Search for files within a directory
- Move and delete files and folders

Let's get started!

Creating Directories and Files

To create a new directory, use the `Path.mkdir()` method. In IDLE's Interactive Window, type the following:

```
>>> from pathlib import Path  
>>> new_dir = Path.home() / "new_directory"  
>>> new_dir.mkdir()
```

After importing the `Path` class, you create a new path to a directory called `new_directory/` in your home folder and assign this path to the `new_dir` variable. Then you use the `.mkdir()` method to create the new directory.

You can now check that the new directory exists and is, in fact, a directory:

```
>>> new_dir.exists()  
True  
  
>>> new_dir.is_dir()  
True
```

If you try to create a directory that already exists, you get an error:

```
>>> new_dir.mkdir()  
Traceback (most recent call last):  
  File "<pysHELL#32>", line 1, in <module>  
    new_dir.mkdir()  
  File "C:\Users\David\AppData\Local\Programs\Python\  
Python38-32\lib\pathlib.py", line 1266, in mkdir  
    self._accessor.mkdir(self, mode)
```

```
FileExistsError: [WinError 183] Cannot create a file when  
that file already exists: 'C:\\\\Users\\\\David\\\\new_directory'
```

When you call the `.mkdir()` method, Python attempts to create the `new_directory/` folder again. Since it already exists, this operation fails and a `FileExistsError` exception is raised.

If you want to create a new directory if it doesn't exist, but avoid raising the `FileExistsError` if it does, then you can set the options `exist_ok` parameter of the `.mkdir()` method to `True`:

```
>>> new_dir.mkdir(exist_ok=True)
```

When you execute `.mkdir()` with the `exist_ok` parameter set to `True`, the directory is created if it does not exist, or nothing happens if it does.

Setting `exist_ok` to `True` when calling `.mkdir()` is equivalent to the following code:

```
>>> if not new_dir.exists():  
...     new_dir.mkdir()
```

Although the above code works just fine, setting the `exist_ok` parameter to `True` is shorter and doesn't sacrifice readability.

Now let's see what happens if you try to create a subdirectory within a directory that does not exist:

```
>>> nested_dir = new_dir / "folder_a" / "folder_b"  
>>> nested_dir.mkdir()  
Traceback (most recent call last):  
  File "<pyshell#38>", line 1, in <module>  
    nested_dir.mkdir()  
  File "C:\\\\Users\\\\David\\\\AppData\\\\Local\\\\Programs\\\\Python\\\\  
Python38-32\\\\lib\\\\pathlib.py", line 1266, in mkdir  
    self._accessor.mkdir(self, mode)  
FileNotFoundException: [WinError 3] The system cannot find the path  

```

The problem is that the directory `folder_a/` does not exist. Typically, to create a directory, all of the parent directories of the target directory `folder_b/` in the path must already exist.

To create any parent directories needed in order to create the target directory, set the optional `parents` parameter of `.mkdir()` to `True`:

```
>>> nested_dir.mkdir(parents=True)
```

Now `.mkdir()` creates the parent directory `folder_a/` so that the target directory `folder_b/` can be created.

By putting all of this together you get the following common pattern for creating directories:

```
path.mkdir(parents=True, exist_ok=True)
```

By setting both the `parents` and `exist_ok` parameters to `True`, the entire path is created, if needed, and no exception is raised if the path already exists.

This pattern is useful, but it may not always be what you want. For example, if the path is input by a user, you may wish to instead catch an exception so that you can ask the user to verify that the path they entered is correct. They might have just mistyped a directory name!

Now let's look at how to create files. Create a new `Path` object called `file_path` for the path `new_directory/file1.txt`:

```
>>> file_path = new_dir / "file1.txt"
```

There is no file in `new_directory/` called `file1.txt`, so the path doesn't exist yet:

```
>>> file_path.exists()  
False
```

You can create the file using the `Path.touch()` method:

```
>>> file_path.touch()
```

This creates a new file called `file1.txt` in the `new_directory/` folder. It doesn't contain any data yet, but the file exists:

```
>>> file_path.exists()  
True  
>>> file_path.is_file()  
True
```

Unlike `.mkdir()`, the `.touch()` method does not raise an exception if the path being created already exists:

```
>>> # Calling .touch() a second time does not raise an exception  
>>> file_path.touch()
```

When you create a file using `.touch()`, the file does not contain any data. You will learn how to write data to a file in Section 11.4: Reading and Writing Files.

You can't create a file in a directory that doesn't exist:

```
>>> file_path = new_dir / "folder_c" / "file2.txt"  
>>> file_path.touch()  
Traceback (most recent call last):  
  File "<pyshell#47>", line 1, in <module>  
    file_path.touch()  
  File "C:\Users\David\AppData\Local\Programs\Python\  
    Python38-32\lib\pathlib.py", line 1256, in touch  
    fd = self._raw_open(flags, mode)  
  File "C:\Users\David\AppData\Local\Programs\Python\  

```

The `FileNotFoundException` exception is raised because the `new_directory/` folder has no `folder_c/` subfolder.

Unlike `.mkdir()`, the `.touch()` method has no `parents` parameter that you can set to automatically create any parent directories. This means that you need to first create any directories needed before calling `.touch()` to create the file.

For instance, you can use `.parent` to get the path to the parent folder for `file2.txt` and then call `.mkdir()` to create the directory:

```
>>> file_path.parent.mkdir()
```

Since `.parent` returns `Path` object, you can chain the `.mkdir()` method to write the entire operation on a single line of code.

With the `folder_c/` directory created, you can successfully create the file:

```
>>> file_path.touch()
```

Now that you know how to create files and directories, let's look at how to get the contents of a directory.

Iterating Over Directory Contents

Using `pathlib`, you can iterate over the contents of a directory. You might need to do this in order to process all of the files in a directory. The word process is vague. It could be reading the file and extracting some data, or compressing files in the directory, or some other operation.

For now, let's focus on how you go about retrieving all of the contents of a directory. You'll learn how to read data from files in the next section.

Everything in a directory is either a file or a subdirectory. The `Path.iterdir()` method returns an iterator over `Path` objects representing each item in the directory.

To use `.iterdir()`, you first need a `Path` representing a directory. Let's use the `new_directory/` folder you created previously in your home di-

rectory and assigned to the new_dir variable:

```
>>> for path in new_dir.iterdir():
...     print(path)
...
C:\Users\David\new_directory\file1.txt
C:\Users\David\new_directory\folder_a
C:\Users\David\new_directory\folder_c
```

Right now, this new_directory/ folder contains three items:

1. A file called file1.txt
2. A directory called folder_c/
3. A directory called folder_a/

Since .iterdir() returns an iterable, you can convert it to a list:

```
>>> list(new_dir.iterdir())
[WindowsPath('C:/Users/David/new_directory/file1.txt'),
 WindowsPath('C:/Users/David/new_directory/folder_a'),
 WindowsPath('C:/Users/David/new_directory/folder_c')]
```

You won't often need to convert this to a list, but we'll do it in subsequent examples to keep the code short. Generally, you'll use .iterdir() in a for loop like you did in the first example.

Notice that .iterdir() only returns items that are directly contained in the new_directory/ folder. That is, you can't see the path to the file that exists in the folder_c/ directory.

There is a way to iterate over the contents a directory and all of its subdirectories, but you can't do it easily with .iterdir(). We'll get to this task in a moment, but first let's talk about how to search for files within a directory.

Searching For Files In a Directory

Sometimes you only need to iterate over files of a certain type, or files with certain naming schemes. You can use the `Path.glob()` method on a path representing a directory to get an iterable over directory contents that meet some criteria.

It might seem strange that a method that searches for files is called `.glob()`. The reason the method is given this name is historical. In early version of the Unix operating system, a program called `glob` was used expand to file path patterns to full file paths.

The `.glob()` method does something similar. You pass to the method a string containing a partial containing a wildcard character and `.glob()` returns a list of file paths that match the pattern.

A **wildcard character** is a special character that acts as a placeholder in a **pattern**. They are replaced with other characters to create a concrete file path. For example, in the pattern `"*.txt"`, the asterisk `*` is a wildcard character that can be replaced with any number of other characters.

The pattern `"*.txt"` matches any file path that ends with `.txt`. That is, if replacing the `*` in the pattern with everything in some file path up to the last four characters results in the original file path, then that file path is a **match** for the pattern `"*.txt"`.

Let's look at an example using the `new_directory/` folder previously assigned to the `new_dir` variable:

```
>>> for path in new_dir.glob("*.txt"):
...     print(path)
...
C:\Users\David\new_directory\file1.txt
```

Like `.iterdir()`, the `.glob()` method returns an iterable of paths, but this time only paths that match the pattern `"*.txt"` are returned. `.glob()` returns only paths that are directly contained in the folder on which it is called.

You can convert the return value of `.glob()` to a list:

```
>>> list(new_dir.glob())
[WindowsPath('C:/Users/David/new_directory/file1.txt')]
```

You will most often use `.glob()` in a `for` loop.

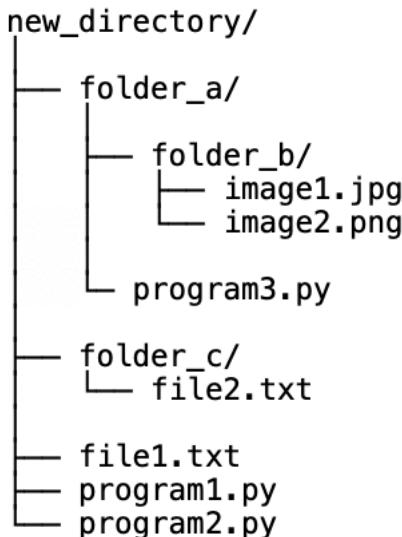
The following table describes some common wildcard characters:

Wildcard Character	Description	Example	Matches	Does Not Match
*	Any number of characters	"*b*"	b, ab, bc, abc	a, c, ac
?	A single character	"?bc"	abc, bbc, cbc	bc, aabc, abcd
[abc]	Matches one character in the brackets	[CB]at	Cat, Bat	at, cat, bat

We'll look at some examples of each of the wildcard characters, but first, let's create a few more files in the `new_directory/` folder so that we have more options to play with:

```
>>> paths = [
...     new_dir / "program1.py",
...     new_dir / "program2.py",
...     new_dir / "folder_a" / "program3.py",
...     new_dir / "folder_a" / "folder_b" / "image1.jpg",
...     new_dir / "folder_a" / "folder_b" / "image2.png",
... ]
>>> for path in paths:
...     path.touch()
...
>>>
```

After executing the above, the `new_directory/` folder has the following structure:



Now that we have a more interesting structure to work with, let's see how `.glob()` works with each of the wildcard characters.

The * Wildcard

The `*` wildcard matches any number of characters in a file path pattern. For example, the pattern "`*.py`" matches all file paths that end in `.py`:

```
>>> list(new_dir.glob("*.py"))
[WindowsPath('C:/Users/David/new_directory/program1.py'),
 WindowsPath('C:/Users/David/new_directory/program2.py')]
```

You can use the `*` wildcard multiple times in a single pattern:

```
>>> list(new_dir.glob("*1*"))
[WindowsPath('C:/Users/David/new_directory/file1.txt'),
 WindowsPath('C:/Users/David/new_directory/program1.py')]
```

The pattern "`*1*`" matches any file path containing the number 1 with any number of characters before and after it. The only files in `new_directory/` that contain the number 1 are `file1.txt` and `program1.py`.

If you leave off the first * from the pattern "*1*" to get the pattern "1*", then nothing gets matched:

```
>>> list(new_dir.glob("1*"))
[]
```

The pattern "1*" matches file paths that start with the number 1 and are followed by any number of characters after it. There are no files in the new_directory/ folder that match this, so .glob() doesn't return anything.

The ? Wildcard

The ? wildcard character matches a single character in a pattern. For example, the pattern "program?.py" will match any file path that starts with the word program followed by a single character and then .py:

```
>>> list(new_dir.glob("program?.py"))
[WindowsPath('C:/Users/David/new_directory/program1.py'),
 WindowsPath('C:/Users/David/new_directory/program2.py')]
```

You can use multiple instances of ? in a single pattern:

```
>>> list(new_dir.glob("?older_?"))
[WindowsPath('C:/Users/David/new_directory/folder_a'),
 WindowsPath('C:/Users/David/new_directory/folder_c')]
```

The pattern "?older_?" matches paths that start with any letter followed by older_ and some other character. In the new_directory/ folder, those paths are the folder_a/ and folder_b/ directories.

You can also combine the * and ? wildcards:

```
>>> list(new_dir.glob("*1.??"))
[WindowsPath('C:/Users/David/new_directory/program1.py')]
```

The pattern "*1.???" matches any file path that contains a 1 followed by a dot (.) and two more characters. The only path in new_directory/ matching this pattern is program1.py. Notice that file1.txt doesn't

match the pattern because the dot is followed by three characters.

The [] Wildcard

The [] wildcard works kind of like the ? wildcard because it matches only a single character. The difference is that instead of matching any single character like ? does, [] only matches characters that are included between the square brackets.

For example, the pattern "program[13].py" matches any path containing the word program, followed by either a 1 or 3 and the extension .py. In the new_directory/ folder, program1.py is the only path matching this pattern:

```
>>> list(new_dir.glob("program[13].py"))
[WindowsPath('C:/Users/David/new_directory/program1.py')]
```

As with the other wildcards, you can use multiple instances of the [] wildcard, as well as combine it with any of the others.

Recursive Matching With The ** Wildcard

The major limitation you've seen with both .iterdir() and .glob() is that they only return paths that are directly contained in the folder on which they are called.

For example, new_dir.glob("*.txt") only returns the file1.txt path in new_directory/. It does not return the file2.txt path in the folder_c/ subdirectory, even though that path matches the "*.txt" pattern.

There is a special wildcard character ** that makes the pattern recursive. The common way to use it is to prefix your pattern with "**/". This tells .glob() to match your pattern in the current directory and any of its subdirectories.

For example, the pattern "**/*.txt" matches both file1.txt and folder_c/file2.txt":

```
>>> list(new_dir.glob("*//*.txt"))
[WindowsPath('C:/Users/David/new_directory/file1.txt'),
 WindowsPath('C:/Users/David/new_directory/folder_c/file2.txt')]
```

Similarly, the pattern "`**/*.py`" matches any `.py` files in `new_directory`/ and any of its subdirectories:

```
>>> list(new_dir.glob("*//*.py"))
[WindowsPath('C:/Users/David/new_directory/program1.py'),
 WindowsPath('C:/Users/David/new_directory/program2.py'),
 WindowsPath('C:/Users/David/new_directory/folder_a/program3.py')]
```

There is also a shorthand method to doing recursive matching called `.rglob()`. To use it, pass the pattern without the `**/` prefix:

```
>>> list(new_dir.rglob("*.py"))
[WindowsPath('C:/Users/David/new_directory/program1.py'),
 WindowsPath('C:/Users/David/new_directory/program2.py'),
 WindowsPath('C:/Users/David/new_directory/folder_a/program3.py')]
```

The `r` in `.rglob()` stands for “recursive.” Some people prefer to use this method instead of prefixing their patterns with `**/` because it is slightly shorter. Both versions are perfectly valid. In this book, we’ll use `.rglob()` instead of the `**/` prefix.

Moving and Deleting Files and Folders

Sometimes you need to move a file or directory to a new location or delete a file or directory all together. You can do this using `pathlib`, but keep in mind that doing so can result in the loss of data, so these operations must be made with extreme care.

To move a file or directory, use the `.replace()` method. For example, the following moves the `file1.txt` file in the `new_directory/` folder to the `folder_a` subfolder:

```
>>> source = new_dir / "file1.txt"
```

```
>>> destination = new_dir / "folder_a" / "file1.txt"
>>> source.replace(destination)
WindowsPath('C:/Users/David/new_directory/folder_a/file1.txt')
```

The `.replace()` method is called on the source path. The destination path is passed to `.replace()` as a single argument. Notice that `.replace()` returns the path to the new location of the file.

Important

If the destination path already exists, `.replace()` overwrites the destination with the source file without raising any kind of exception. This can cause undesired loss of data if you aren't careful.

You may want to first check if the destination file exists, and move the file only in the case that it does not:

```
if not destination.exists():
    source.replace(destination)
```

You can also use `.replace()` to move or rename an entire directory. For instance, the following renames the `folder_c` subdirectory of `new_directory` to `folder_d`:

```
>>> source = new_dir / "folder_c"
>>> destination = new_dir / "folder_d"
>>> source.replace(destination)
WindowsPath('C:/Users/David/new_directory/folder_d')
```

Again, if the destination folder already exists, it is completely replaces with the source folder, which could result in the loss of quite a bit of data.

To delete a file, use the `.unlink()` method:

```
>>> file_path = new_dir / "program1.py"
>>> file_path.unlink()
```

This deletes the `program1.py` file in the `new_directory/` folder, which you can check with `.exists()`:

```
>>> file_path.exists()  
False
```

You can also see it removed with `.iterdir()`:

```
>>> list(new_dir.iterdir())  
[WindowsPath('C:/Users/David/new_directory/folder_a'),  
 WindowsPath('C:/Users/David/new_directory/folder_d'),  
 WindowsPath('C:/Users/David/new_directory/program2.py')]
```

If the path that you call `.unlink()` does not exists, a `FileNotFoundError` exception is raised:

```
>>> file_path.unlink()  
Traceback (most recent call last):  
  File "<pyshell#94>", line 1, in <module>  
    file_path.unlink()  
  File "C:\Users\David\AppData\Local\Programs\Python\  
    Python38-32\lib\pathlib.py", line 1303, in unlink  
    self._accessor.unlink(self)  
FileNotFoundException: [WinError 2] The system cannot find the file  
specified: 'C:\\\\Users\\\\David\\\\new_directory\\\\program1.py'
```

If you want to ignore the exception, set the optional `missing_ok` parameter to `True`:

```
>>> file_path.unlink(missing_ok=True)
```

In this case, nothing actually happens because the file located at `file_path` does not exist.

Important

When you delete a file it is gone forever. Make sure you really want to delete it before you proceed!

.unlink() only works for paths representing files. To remove a directory, use the .rmdir() method. Keep in mind that the folder must be empty, otherwise an OSError exception is raised:

```
>>> folder_d = new_dir / "folder_d"  
>>> folder_d.rmdir()  
Traceback (most recent call last):  
  File "<pyshell#97>", line 1, in <module>  
    folder_d.rmdir()  
  File "C:\Users\David\AppData\Local\Programs\Python\  
    Python38-32\lib\pathlib.py", line 1314, in rmdir  
    self._accessor.rmdir(self)  
OSError: [WinError 145] The directory is not empty:  

```

In the case of folder_d/, it only contains a single file called file2.txt. To delete folder_d/, first delete all of the files it contains:

```
>>> for path in folder_d.iterdir():  
...     path.unlink()  
...  
>>> folder_d.rmdir()
```

Now folder_d/ is deleted:

```
>>> folder_d.exists()  
False
```

If you need to delete an entire directory, even if it is non-empty, then pathlib won't help you much. However, you can use the rmtree() function from the built-in shutil module:

```
>>> import shutil  
>>> folder_a = new_dir / "folder_a"  
>>> shutil.rmtree(folder_a)
```

Recall that folder_a/ contains a subfolder folder_b/ which itself contains two files called image1.jpg and image2.png.

When you pass the `folder_a` path object to `rmtree()`, the `folder_a/` and all of its contents are deleted:

```
>>> # The folder_a/ directory no longer exists
>>> folder_a.exists()
False

>>> # Searching for `image*.*` files returns nothing
>>> list(new_dir.rglob("image*.*"))
[]
```

In this section you covered quite a bit of ground. You learned how to do several common file system operations, such as:

- Creating files and directories
- Iterating over the contents of a directory
- Searching for files and folders using wildcards
- Moving and deleting files and folders

All of these are common tasks. It is extremely important, however, to remember that your programs are guests on another persons computer. If you aren't careful, you can inadvertently cause damage to a user's computer resulting the loss of important documents and other data.

When working with the file system you should always use caution. When in doubt, check that file paths exist or do not exists before performing some operation, and always check with the user that what you are about to do is OK!

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Create a new directory in your home folder called `my_folder/`.
2. Inside `my_folder/` create three files:

- file1.txt
 - file2.txt
 - image1.png
3. Move the file image1.png to a new directory called images/ inside of the my_folder/ directory.
 4. Delete the file file1.txt
 5. Delete the my_folder/ directory.

[Leave feedback on this section »](#)

12.4 Challenge: Move All Image Files To a New Directory

In the Chapter 12 Practice Files folder, there is a subfolder called documents/. The directory contains several files and subfolders. Some of the files are images ending with either the .png, .gif, or the .jpg file extension.

Create a new folder in the Practice Files folder called images/ and move all image files to that folder. When you are done, the new folder should have four files in it:

1. image1.png
2. image2.gif
3. image3.png
4. image4.jpg

You can find the solutions to this code challenge and many other bonus resources online at realpython.com/python-basics/resources.

[Leave feedback on this section »](#)

12.5 Reading and Writing Files

Files are abundant in the modern world. They are the medium through which data is digitally stored and transferred. Chances are, you've opened dozens, if not hundreds, of files just today.

In this section, you'll learn how to read and write files with Python.

What Is a File?

A **file** is a sequence of bytes and a **byte** is a number between 0 and 255. That is, a file is a sequence of integer values.

The bytes in a file must be **decoded** into something meaningful in order to understand the contents of the file.

Python has standard library modules for working with text, csv, and audio files. There are a number of third-party packages available for working with other file types.

You'll learn how to install third-party packages in Chapter 13: *Installing Packages with Pip*. You'll also see how to work with PDF files in Chapter 14: *Creating and Modifying PDF Files*.

In this section, you'll learn how to work with plain text files.

Understanding Text Files

Text files are files that contain only text. They are perhaps the easiest files to work with. There are two issues, though, that can be frustrating when working with text files:

1. Character encoding
2. Line endings

Before jumping into reading and writing text files, let's look at what these issues are so that you know how to deal with them effectively.

Character Encoding

Text files are stored on disk as a sequence of bytes. Each byte, or group of bytes in some cases, represents a different character in the file.

When text files are written, characters typed on the keyboard are converted to bytes in a process called **encoding**. When a text file is read, the bytes are **decoded** back into text.

The integer a character is associated to is determined by the file's **character encoding**. There are many character encodings. Four of the most widely used character encodings are:

1. ASCII
2. UTF-8
3. UTF-16
4. UTF-32

Some character encodings, such as ASCII and UTF-8, encode characters the same way. For example, numbers and English letters are encoded the same way in both ASCII and UTF-8.

The difference between ASCII and UTF-8 is that UTF-8 can encode more characters than ASCII. ASCII can't encode characters like ñ or ü, but UTF-8 can. This means you can decode ASCII encoded text with UTF-8, but you can't always decode UTF-8 encoded text with ASCII.

Important

Serious problems may occur when different encodings are used to encode and decode text.

For instance, text encoded as UTF-8 that is decoded with UTF-16 may be interpreted as an entirely different language than originally intended!

For a thorough introduction to character encodings, check out Real Python's [Unicode & Character Encodings in Python: A Painless Guide](#).

Knowing what encoding a file uses is important, but it isn't always obvious. On modern Windows computers, txt files are usually encoded with UTF-16 or UTF-8. On macOS and Ubuntu Linux, the default character encoding is usually UTF-8.

For the remainder of this section, we'll assume that the character encoding of all text files that we work with is UTF-8. If you encounter problems, you may need to alter the examples to use a different encoding.

Line Endings

Each line in a text file ends with one or two characters that indicate the line has ended. These characters aren't usually displayed in a text editor, but they exist as bytes in the file data.

The two characters used to represent line endings are the **carriage return** and **line feed** characters. In Python strings, these characters are represented by the escape sequence `\r` and `\n`, respectively.

On Windows, line endings are represented by default with both a carriage return and a line feed. On macOS and most Linux distributions, line endings are represented with just a single line feed character.

When you read a Windows file on macOS or Linux you will sometimes see extra blank lines between lines of text. This is because the carriage

return also represents a line ending on macOS and Linux.

For example, suppose the following text file was created in Windows:

```
Pug\r\n
Jack Russell Terrier\r\n
English Springer Spaniel\r\n
German Shepherd\r\n
```

On macOS or Ubuntu, this file is interpreted with double spacing between lines:

```
Pug\r
\n
Jack Russell Terrier\r
\n
English Springer Spaniel\r
\n
German Shepherd\r
\n
```

In practice, the differences between line endings on different operating systems is rarely problematic. Python can handle line ending conversions for you automatically, so you don't have to worry about it too often.

Python File Objects

Files are represented in Python with **file objects**, which are instances of classes designed to work with different types of files.

Python has a couple of different types of file objects:

1. Text file objects are used for interacting with text files
2. Binary file objects are used for working directly with the bytes contained in files

Text file objects handle encoding and decoding bytes for you. All you

need to do is specify which character encoding to use. On the other hand, binary file objects do not perform any kind of encoding or decoding.

There are two ways to create a file object in Python:

1. The `Path.open()` method
2. The `open()` built-in function

Let's look at each of these.

The `Path.open()` Method

To use the `Path.open()` method, you first need a `Path` object. In IDLE's interactive window, execute the following:

```
>>> from pathlib import Path  
>>> path = Path.home() / "hello.txt"  
>>> path.touch()  
>>> file = path.open(mode="r", encoding="utf-8")
```

First, a `Path` object for the `hello.txt` file is created and assigned to the `path` variable. Then `path.touch()` creates the file in your home directory. Finally, `.open()` returns a new file object representing the `hello.txt` file and assigns it to the `file` variable.

Two keyword parameters used to open the file:

1. The `mode` parameter determines in which mode the file should be opened. The "`r`" argument opens the file in **read mode**.
2. The `encoding` parameter determines the character encoding used to decode the file. The argument "`utf-8`" represents the UTF-8 character encoding.

You can inspect the `file` variable to see that it is assigned to a text file object:

```
>>> file
<_io.TextIOWrapper name='C:\Users\David\hello.txt' mode='r'
encoding='utf-8'>
```

Text file objects are instances of the `TextIOWrapper` class. You will never need to instantiate this class directly, since you can create them with the `Path.open()` method.

There are a number of different modes you can use to open a file. These are described in the following table:

Mode	Description
"r"	Creates a text file object for reading and raises an error if the file can't be opened.
"w"	Creates a text file object for writing and overwrites all existing data in the file.
"a"	Creates a text file object for appending data to the end of a file.
"rb"	Creates a binary file object for reading and raises an error if the file can't be opened.
"wb"	Creates a binary file object for writing and overwrites all existing data in the file.
"ab"	Creates a binary file object for appending data to the end of the file

The strings for some of the most commonly used character encodings can be found in the table below:

String	Character Encoding
"ascii"	ASCII
"utf-8"	UTF-8
"utf-16"	UTF-16
"utf-32"	UTF-32

When you create a file object with `.open()`, Python maintains a link to the file resource until you either explicitly tell Python to close the file, or the program ends.

Important

You should always explicitly tell Python to close a file.

Forgetting to close opened files like littering. When your program stops running, it shouldn't leave unnecessary waste laying around the system.

To close a file, use the file object's `.close()` method:

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

Using `Path.open()` is the preferred way to open a file when you have an existing `Path` object, but there is also a built-in function called `open()` that you can use to open a file.

The `open()` Built-in

The built-in `open()` function works almost exactly like the `Path.open()` method, except that it's first parameter is a string containing the path the file you want to open.

First, create a new variable called `file_path` and assign to it a string containing the path to the `hello.txt` file you created above:

```
>>> file_path = "C:/Users/David/hello.txt"
```

Note that you'll need to change the path to match the path of the file on your own computer.

Next, create a new file object using the `open()` built-in and assign it to the variable `file`:

```
>>> file = open(file_path, mode="r", encoding="utf-8")
```

The first parameter of `open()` must be a path string. The `mode` and `encoding` parameters are the same as the parameters for the `Path.open()` method. In this example, `mode` is set to "r" for read mode, and `encoding` is set to "utf-8".

Just like the file object returned by `Path.open()`, the file object returned by `open()` is a `TextIOWrapper` instance:

```
>>> file
<_io.TextIOWrapper name='C:/Users/David/hello.txt' mode='r' encoding='utf-8'>
```

To close the file, use the file object's `.close()` method:

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

For the most part, you'll use the `Path.open()` method to open a file from an existing `pathlib.Path` object. However, if you don't need all of the functionality of the `pathlib` module, then `open()` is a great way to quickly create a file object.

The `with` Statement

When you open a file, your program is accessing data external to the program itself. The operating system must manage the connection between your program and physical file itself. When you call a file object's `.close()` method, the operating system knows to close the connection.

If your program crashes between the time that a file is opened and when it is closed, the system resources maintained by the connection may continue to live on until the operating system realizes that it's no longer needed.

To ensure that file system resources are cleaned up even if a program crashes, you can open a file in a `with` statement. The pattern for using the `with` statement looks like this:

```
with path.open(mode="r", encoding="utf-8") as file:  
    # Do something with file
```

The `with` statement has two parts: a header and a body. The header always starts with the `with` keyword and ends with a colon (:). The return value of `path.open()` is assigned to the variable name after the `as` keyword.

After the `with` statement header is an indented block of code. When code execution leaves the indented block, the file object assigned to `file` is closed automatically, even if any exception is raised during execution of code inside of the block.

`with` statements also work with the `open()` built-in:

```
with open(file_path, mode="r", encoding="utf-8") as file:  
    # Do something with file
```

There really is no reason *not* to open files in a `with` statement. It is considered the Pythonic way for working with files. For the rest of this book, we will use this pattern whenever opening a file.

Reading Data From a File

Using a text editor, open the `hello.txt` file in your home directory that you previously created and type the text `Hello World` into it. Then save the file.

In IDLE's interactive window, type the following:

```
>>> path = Path.home() / "hello.txt"  
>>> with path.open(mode="r", encoding="utf-8") as file:  
...     text = file.read()  
...  
>>>
```

The file object created by `path.open()` is assigned to the `file` variable. Inside of the `with` block, the file object's `.read()` method reads the text from the file and assigns the result to the variable `text`.

The value returned by `.read()` is a string object with the value "Hello World":

```
>>> type(text)
<class 'str'>
>>> text
'Hello World'
```

The `.read()` method reads all of the text in the file and returns it as a string value.

If there are multiple lines of text in the file, each line in the string is separated with a newline character `\n`. In a text editor, open the `hello.txt` file again and put the text "Hello again" on the second line. Then save the file.

Back in IDLE's interactive window, read the text from the file again:

```
>>> with path.open(mode="r", encoding="utf-8") as file:
...     text = file.read()
...
>>> text
'Hello World\nHello again'
```

The text from each line has a `\n` character in between.

Instead of reading the entire file at once, you can process each line of the file one at a time:

```
>>> with path.open(mode="r", encoding="utf-8") as file:
...     for line in file.readlines():
...         print(line)
...
Hello World

Hello again
```

The `.readlines()` method returns an iterable of lines from the file. At each step of the `for` loop the next line of text in the file is returned and

printed.

Notice that an extra blank line is printed between the two lines of text. This doesn't have anything to do with line endings in the file. It happens because the `print()` function automatically inserts a newline character at the end of every string it prints.

To print the two lines without the extra blank line, set the `print()` function's optional `end` parameter to an empty string:

```
>>> with path.open(mode="r", encoding="utf-8") as file:  
...     for line in file.readlines():  
...         print(line, end="")  
...  
Hello World  
Hello again
```

There are many times you might want to use `.readlines()` instead of `.read()`. For example, each line in a file might represent a single record. You can loop over the lines of text in the file with `.readlines()` and process them as needed.

If you try to read from a file that does not exist, both `.open()` and `open()` raise a `FileNotFoundException`:

```
>>> path = Path.home() / "new_file.txt"  
>>> with path.open(mode="r", encoding="utf-8") as file:  
    text = file.read()  
  
Traceback (most recent call last):  
  File "<pyshell#197>", line 1, in <module>  
    with path.open(mode="r", encoding="utf-8") as file:  
  File "C:\Users\David\AppData\Local\Programs\Python\  
Python38-32\lib\pathlib.py", line 1200, in open  
    return io.open(self, mode, buffering, encoding, errors, newline,  
File "C:\Users\David\AppData\Local\Programs\Python\
```

```
Python38-32\lib\pathlib.py", line 1054, in _opener
    return self._accessor.open(self, flags, mode)
FileNotFoundError: [Errno 2] No such file or directory:
'C:\\\\Users\\\\David\\\\new_file.txt'
```

Next, let's see how to write data to a file.

Writing Data To a File

To write data to a plain text file, you pass a string to a file object's `.write()` method. The file object must be opened in **write mode** by passing the value "w" to the `mode` parameter.

For instance, the following writes the text "Hi there!" to the `hello.txt` file in your home directory:

```
>>> with path.open(mode="w", encoding="utf-8") as file:
...     file.write("Hi there!")
...
9
>>>
```

Notice that the integer 9 is displayed after executing the `with` block. That's because `.write()` returns the numbers of characters that are written. The string "Hi there!" has nine characters, so `.write()` returns the number 9.

When the text "Hi there!" is written to the `hello.txt` file, any existing contents are written over. It's as if you deleted the old `hello.txt` file and created a new one.

Important

When you set `mode="w"` in `.open()`, the contents of the original file are overwritten. This results in the loss of all of the original data in the file!

You can verify that the file only contains the text "Hi there!" by reading and displaying the contents of the file:

```
>>> with path.open(mode="r", encoding="utf-8") as file:  
...     text = file.read()  
...  
>>> print(text)  
Hi there!
```

You can **append** data to the end of a file by opening the file in **append mode**:

```
>>> with path.open(mode="a", encoding="utf-8") as file:  
...     file.write("\nHello")  
...  
6
```

When a file is opened in append mode new data is written to the end of the file and old data is left intact. The newline character is put at the beginning of the string so that the word "Hello" is printed on a new line at the end of the file.

Without a newline character at the beginning of the string, the word "Hello" would be printed on the same line as any existing text at the end of the file.

You can check that the word "Hello" is written to the second line by opening and reading from the file:

```
>>> with path.open(mode="r", encoding="utf-8") as file:  
...     text = file.read()  
...  
>>> print(text)  
Hi there!  
Hello
```

You can write multiple lines to a file at the same time using the `.writelines()` method. First, create a list of strings:

```
>>> lines_of_text = [  
...     "Hello from Line 1\n",
```

```
...     "Hello from Line 2\n",
...     "Hello from Line 3 \n",
... ]
```

Then open the file in write mode and use the `.writelines()` method to write each string in the list to the file:

```
>>> with path.open(mode="w", encoding="utf-8") as file:
...     file.writelines(lines_of_text)
...
>>>
```

Each string in the `lines_of_text` list is written to the file. Notice that each string ends with the newline character (`\n`). That's because `.writelines()` doesn't automatically insert each string in the list on a new line.

If you open a non-existent path in write mode, Python attempts to automatically create the file. If all of the parent folders in the path exist, then the file can be created without problem:

```
>>> path = Path.home() / "new_file.txt"
>>> with path.open(mode="w", encoding="utf-8") as file:
...     file.write("Hello!")
...
6
```

Since the `Path.home()` directory exists, a new file called `new_file.txt` is created automatically.

However, if one of the parent directories does not exist, then `.open()` will raise a `NotFoundError`:

```
>>> path = Path.home() / "new_folder" / "new_file.txt"
>>> with path.open(mode="w", encoding="utf-8") as file:
...     file.write("Hello!")
...
Traceback (most recent call last):
```

```
File "<pyshell#172>", line 1, in <module>
    with path.open(mode="w", encoding="utf-8") as file:
File "C:\Users\David\AppData\Local\Programs\Python\
Python38-32\lib\pathlib.py", line 1200, in open
    return io.open(self, mode, buffering, encoding, errors, newline,
File "C:\Users\David\AppData\Local\Programs\Python\
Python38-32\lib\pathlib.py", line 1054, in _opener
    return self._accessor.open(self, flags, mode)
 FileNotFoundError: [Errno 2] No such file or directory:
'C:\\\\Users\\\\David\\\\new_folder\\\\new_file.txt'
```

If you want to write to a path with parent folders that may not exist, call the `.mkdir()` method with the `parents` parameter set to `True` before opening the file in write mode:

```
>>> path.parent.mkdir(parents=True)
>>> with path.open(mode="w", encoding="utf-8") as file:
...     file.write("Hello!")
...
6
```

In this section you covered a lot of ground. You learned that all files are sequences of bytes, which are integers with values between 0 and 255.

You also learned about character encodings, which are used to translate between bytes and text, and differences between line endings on different operating systems.

Finally, you saw how to read and write text files using the `Path.open()` method and the `open()` built-in.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Write the following text to file called `starships.txt` in your home

directory:

```
Discovery  
Enterprise  
Defiant  
Voyager
```

Each word should be on a separate line.

2. Read the file `starships.txt` you created in Exercise 1 and print each line of text in the file. The output should not have extra blank lines between each word.
3. Read the file `startships.txt` and print the names of the starships that start with the letter D.

[Leave feedback on this section »](#)

12.6 Read and Write CSV Data

Suppose you had a temperature sensor in your house that records the temperature every four hours. Over the course of a day, six temperature readings are taken.

You can store each temperature reading in a list:

```
>>> temperature_readings = [68, 65, 68, 70, 74, 72]
```

Each day a new list of numbers is generated. To store these values to a file, you can write the values from each day on a new line in a text file and separate each value with a comma.

```
>>> from pathlib import Path  
>>> file_path = Path.home() / "temperatures.txt"  
>>> with file_path.open(mode="a", encoding="utf-8") as file:  
...     file.write(str(temperature_readings[0]))  
...     for temp in temperature_readings[1:]:  
...         file.write(f",{temp}")  
...
```

```
2
3
3
3
3
3
3
>>>
```

This creates a file called `temperatures.csv` in your home directory and opens it in append mode. On a new line at the end of the file, the first value in the `temperature_readings` list is written to the file. Then each remaining value in the list is written, preceded by a comma, to the same line.

The final string of text written to the file "68,65,68,70,74,72". You can verify this by reading the text:

```
>>> with file_path.open(mode="r", encoding="utf-8") as file:
...     text = file.read()
...
>>> text
'68,65,68,70,74,72'
```

The format that in which you have saved the values is called **comma-separated value**, or **CSV** for short. The `temperatures.csv` file is called a csv file.

csv files are a great way to store records of sequential data because you can recover each row of the csv value as a list:

```
>>> temperatures = text.split(",")
>>> temperatures
['68', '65', '68', '70', '74', '72']
```

In Section 9.2: *Lists Are Mutable Sequences*, you learned how to create a list from a string using the `.split()` string method. In the example above, a new list is created from the `text` read from the `temperatures.csv` file.

The values in the `temperatures` list are strings, not integers like the values originally written to the file. This is because values read from a text file are always read as strings.

You can change convert the strings to integers using a list comprehension:

```
>>> int_temperatures = [int(temp) for temp in temperatures]
>>> int_temperatures
[68, 65, 68, 70, 74, 72]
```

You have now recovered the list that you originally wrote to the `temperatures.csv` file!

What these examples illustrate is that a csv file is a plain text file. Using techniques from (Section 11.4: Reading and Writing Files)[#read-write-files], you can store sequences of values in the rows of the csv file, and then read from the file to recover the data.

Reading and writing csv files is so common that the Python standard library has a module called `csv` to lessen the workload required for working with csv files. In the following sections, you'll learn how to use the `csv` module to write to and read from csv files.

The `csv` Module

The `csv` module can be used to read and write csv files. We'll re-work the previous example using the `csv` module so that you see how it works and what operations it takes care of for you.

To get started, import the `csv` module in IDLE's interactive window:

```
>>> import csv
```

Let's start by creating a new csv file containing several days worth of temperature data.

Writing csv Files With `csv.writer`

Create a list of lists containing temperature readings from three days:

```
>>> daily_temperatures = [
...     [68, 65, 68, 70, 74, 72],
...     [67, 67, 70, 72, 72, 70],
...     [68, 70, 74, 76, 74, 73],
... ]
```

Now open the `temperatures.csv` file in write mode:

```
>>> file_path = Path.home() / "temperatures.csv"
>>> file = file_path.open(mode="w", encoding="utf-8")
```

Instead of using a `with` statement, a file object is created and assigned to the `file` variable so that we can inspect each step of the writing process as we go.

Now create a new csv writer object by passing the file object `file` to `csv.writer()`:

```
>>> writer = csv.writer()
```

`csv.writer()` return a csv writer object with methods for writing data to the csv file.

For instance, you can use the `writer.writerow()` method write a list to a new row in the csv file:

```
>>> for temp_list in daily_temperatures:
...     writer.writerow(temp_list)
...
19
19
19
```

Just like a file object's `.write()` method, `.writerow()` returns the number of characters written to the file. Each list in `daily_temperatures` gets

converted to a string containing the temperatures separated by commas, and each of these strings has 19 characters.

Now close the file:

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

If you open the `temperatures.csv` file in a text editor, you will see the following text in the file:

```
68,65,68,70,74,72  
67,67,70,72,72,70  
68,70,74,76,74,73
```

In the examples above, you did not use a `with` statement to write to the file so that you can inspect each operation in IDLE's interactive window. You won't typically do this in practice, so here's what the code looks like using the `with` statement:

```
with file_path.open(mode="w", encoding="utf-8") as file:  
    writer = csv.writer(file)  
    for temp_list in daily_temperatures:  
        writer.writerow(temp_list)
```

The main advantage of using `csv.writer` to write to a csv file is that you don't need to worry about converting values to strings before writing them to the file. The `csv.writer` object handles this for you, which results in shorter and cleaner code.

`.writerow()` writes a single row to the csv file, but you can write multiple rows at one using the `.writerows()` method. This shortens the code even more when your data is already in a list of lists:

```
with file_path.open(mode="w", encoding="utf-8") as file:  
    writer = csv.writer(file)  
    writer.writerows(daily_temperatures)
```

Now let's read from the `temperatures.csv` file to recover the `daily_temperatures` list of lists that was used to create the file.

Reading csv Files With `csv.reader`

To read a csv file with the `csv` module, use the `csv.reader` class. Like `csv.writer` objects, `csv.reader` objects are instantiated from a file object:

```
>>> file = file_path.open(mode="r", encoding="utf-8")
>>> reader = csv.reader(file)
```

`csv.reader()` returns a csv reader object that can be used to iterate over the rows of the csv file:

```
>>> for row in reader:
...     print(row)
...
['68', '65', '68', '70', '74', '72']
['67', '67', '70', '72', '72', '70']
['68', '70', '74', '76', '74', '73']
>>> file.close()
```

Each row of the csv file is returned as a list of strings. To recover the `daily_temperatures` list of lists, you'll need to convert each list of strings to a list of integers using a list comprehension.

Here's a full example using that open the csv file in a `with` statement, reads each row in the csv file, converts the list of strings to a list of integers, and stores each list of integers in a list of lists called `daily_temperatures`:

```
>>> # Create an empty list
>>> daily_temperatures = []
>>> with file_path.open(mode="r", encoding="utf-8") as file:
...     reader = csv.reader(file)
...     for row in reader:
...         # Convert row to list of integers
...         int_row = [int(value) for value in row]
...         # Append the list of integers to daily_temperatures list
...         daily_temperatures.append(int_row)
```

```
...
>>> daily_temperatures
[[68, 65, 68, 70, 74, 72], [67, 67, 70, 72, 72, 70],
 [68, 70, 74, 76, 74, 73]]
```

It is much easier to work with csv files using the `csv` module than it is using the standard tools for reading and writing plain text files.

Sometimes, though, csv files are more complex than a file with rows of values that all have the same type. Each row may represent a record with various fields, and the first row in the file may be a **header row** with the names of the fields.

Reading and Writing CSV Files With Headers

Here's an example of a csv file with a header row containing multiple data types:

```
name,department,salary
Lee,Operations,75000.00
Jane,Engineering,85000.00
Diego,Sales,80000.00
```

The first line of the file contains field names. Each following line contains a record with a value for each field.

It's possible to read csv files such as the one above using `csv.reader()`, but you have to keep track of the header row, and each row is returned as a list without the field names attached to it. It makes more sense to return each row as a dictionary whose keys are the field names and values are the field values in the row. This is precisely what `csv.DictReader` objects do!

Using a text editor, create a new csv file called `employees.csv` and save the text from the example csv file above to it. Save the file to your computer's home directory.

In IDLE's interactive window, open the `employees.csv` file and create a new `csv.DictReader` object:

```
>>> file_path = Path.home() / "employees.csv"
>>> file = file_path.open(mode="r", encoding="utf-8")
>>> reader = csv.DictReader(file)
```

When you create a `DictReader` object, the first row of the csv file is assumed to contain the field names. These values get stored in a list and assigned to the `DictReader` instance's `.fieldnames` attribute:

```
>>> reader.fieldnames
['name', 'department', 'salary']
```

Just like `csv.reader` objects, `DictReader` objects are iterable:

```
>>> for row in reader:
...     print(row)
...
{'name': 'Lee', 'department': 'Operations', 'salary': '75000.000'}
{'name': 'Jane', 'department': 'Engineering', 'salary': '85000.00'}
{'name': 'Diego', 'department': 'Sales', 'salary': '80000.00'}
>>> file.close()
```

Instead of returning each row as a list, `DictReader` objects return each row as a dictionary. The dictionary's keys are the field names, and the values are the field values from each row in the csv file.

Notice that the `salary` field gets read as a string. Since csv files are plain text files, the values are always read as strings. You'll need to convert the strings to different data types as needed.

For example, you can process each row with a function that converts keys to the correct data types:

```
>>> def process_row(row):
...     row["salary"] = float(row["salary"])
...     return row
...
>>> with file_path.open(mode="r", encoding="utf-8") as file:
```

```
...     reader = csv.DictReader(file)
...     for row in reader:
...         print(process_row(row))
...
{'name': 'Lee', 'department': 'Operations', 'salary': 75000.0}
{'name': 'Jane', 'department': 'Engineering', 'salary': 85000.0}
{'name': 'Diego', 'department': 'Sales', 'salary': 80000.0}
```

The `process_row()` function takes a row dictionary read from the csv file and returns a new dictionary with the "salary" key converted to a floating point number.

You can write csv files with headers using the `csv.DictWriter` class, which writes dictionaries with shared keys to rows in a csv file.

The following list of dictionaries represents a small database of people and their ages:

```
>>> people = [
...     {"name": "Veronica", "age": 29},
...     {"name": "Audrey", "age": 32},
...     {"name": "Sam", "age": 24},
... ]
```

To store the data in the `people` list to a csv file, open a new file called `people.csv` in write mode and create a new `csv.DictWriter` object from the file object:

```
>>> file_path = Path.home() / "people.csv"
>>> file = file_path.open(mode="w", encoding="utf-8")
>>> writer = csv.DictWriter(file, fieldnames=["name", "age"])
```

When you instantiate a new `DictWriter` object, the first parameter is the file object for writing the csv data. The `fieldnames` parameter, which is required, is a list of strings of the field names.

Note

In the example above, the string literal `["name", "age"]` is passed to the `fieldnames` parameter, but you don't have to use a string literal. For example, you could also set the `fieldnames` parameter to `people[0].keys()`.

This can be useful if the field names are not known when writing the program, or if there are so many fields that a list literal is impractical.

Just like `csv.writer` objects, `DictReader` objects have a `.writerow()` method for writing a single row to the csv file and a `.writerows()` method for writing several rows at once. But `DictReader` objects have a third method called `.writeheader()` that writes the header row to the csv file:

```
>>> writer.writeheader()
```

```
10
```

`.writeheader()` returns the number of characters written to the file, which is 10 in this case. Writing the header row is optional, but is recommended because it helps define what the data contained in the csv file represents. It also makes it easy to read the rows from the csv file as dictionaries using the `DictReader` class.

With the header written, you can write the data in the `people` list to the csv file using `.writerows()`:

```
>>> writer.writerows(people)
>>> file.close()
```

You now have a file in your home directory called `people.csv` containing the following data:

```
name,age
Veronica,29
Audrey,32
```

Sam,24

csv files are a flexible and convenient way of storing data. They are used frequently in business worldwide, and knowing how to work with them is a valuable skill!

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Write a script that writes the following list of lists to a file called `numbers.csv` in your home directory:

```
numbers = [  
    [1, 2, 3, 4, 5],  
    [6, 7, 8, 9, 10],  
    [11, 12, 13, 14, 15],  
]
```

2. Write a script that reads the numbers in the `numbers.csv` file from Exercise 1 into a list of lists of integers called `numbers`. Print the list of lists. Your output should like the following:

```
[[1, 2, 3, 4, 5], [6, 7, 8, 9, 10], [11, 12, 13, 14, 15]]
```

3. Write a script that writes the following list of dictionaries to a file called `favorite_colors.csv` in your home directory:

```
favorite_colors = [  
    {"name": "Joe", "favorite_color": "blue"},  
    {"name": "Anne", "favorite_color": "green"},  
    {"name": "Bailey", "favorite_color": "red"},  
]
```

The output csv file should have the following format:

```
name,favorite color  
Joe,blue  
Anne,green  
Bailey,red
```

4. Write a script that reads the data from the `favorite_colors.csv` file from Exercise 3 into a list of dictionaries called `favorite_colors`.

Print the list of dictionaries. The output should look something like this:

```
[{"name": "Joe", "favorite_color": "blue"},  
 {"name": "Anne", "favorite_color": "green"},  
 {"name": "Bailey", "favorite_color": "red"}]
```

[Leave feedback on this section »](#)

12.7 Challenge: Create a High Scores List

In the Chapter 12 Practice Files folder, there is a csv file called `scores.csv` containing data about game players and their scores. The first few lines of the file look like this:

```
name, score  
LLCoolDave, 23  
LLCoolDave, 27  
red, 12  
LLCoolDave, 26  
tom123, 26
```

Write a script that reads the data from this csv file and creates a new file called `high_scores.csv` where each row contains the player name and their highest score.

The output csv file should look like this:

```
name, high_score  
LLCoolDave, 27  
red, 12  
tom123, 26  
0_0, 7  
Misha46, 25  
Empiro, 23  
MaxxT, 25
```

L33tH4x, 42
johnsmith, 30

You can find the solutions to this code challenge and many other bonus resources online at realpython.com/python-basics/resources.

[Leave feedback on this section »](#)

12.8 Summary and Additional Resources

In this chapter you learned about the file system and file paths and how to work with them using the Python standard library’s `pathlib` module. You saw how to create new `Path` objects, access path components, and how to create, move, and delete files and folders.

You also learned how to read and write plain text files using the `Path.open()` method and `open()` built-in, and how to work with comma-separated value, or csv, files using the Python standard library `csv` module.

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-11

Additional Resources

To get even more practice working with files, check out these resources:

- [Reading and Writing Files in Python \(Guide\)](#)
- [Working With Files in Python](#)

- Recommended resources on realpython.com

[Leave feedback on this section »](#)

Chapter 13

Installing Packages With Pip

Up to this point, you have been working within the bounds of the **Python standard library**. In the remaining half of this course, you will work with various **packages** that are not included with Python by default.

Many programming languages offer a **package manager** that automates the process of installing, upgrading, and removing third-party packages. Python is no exception.

The *de facto* package manager for Python is called `pip`. Historically, `pip` had to be downloaded and installed separately from Python. As of Python 3.4, it is now included with most distributions of the language.

In this chapter, you will learn:

- How to install and manage third-party packages with `pip`
- What the benefits and risks of third-party packages are

Let's go!

[Leave feedback on this section »](#)

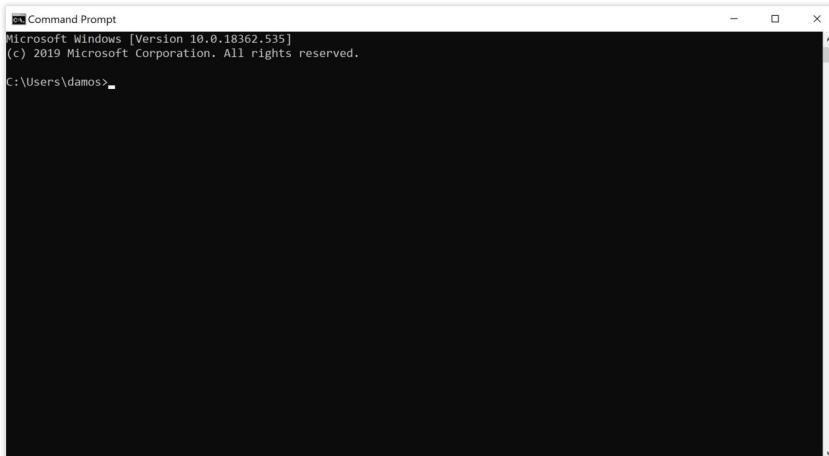
13.1 Installing Third-Party Packages With Pip

Python's package manager `pip` is used to install and manage third party packages. It is a separate program from Python, although it's likely that `pip` was installed on your computer whenever you downloaded and installed Python.

`pip` is a **command line tool**. That means you must run it from a command line or terminal program. How you open a terminal program depends on your operating system.

Windows

Press the Windows key and type `cmd` and press `Enter` to open the `Command Prompt` application. This opens a window that looks like this:

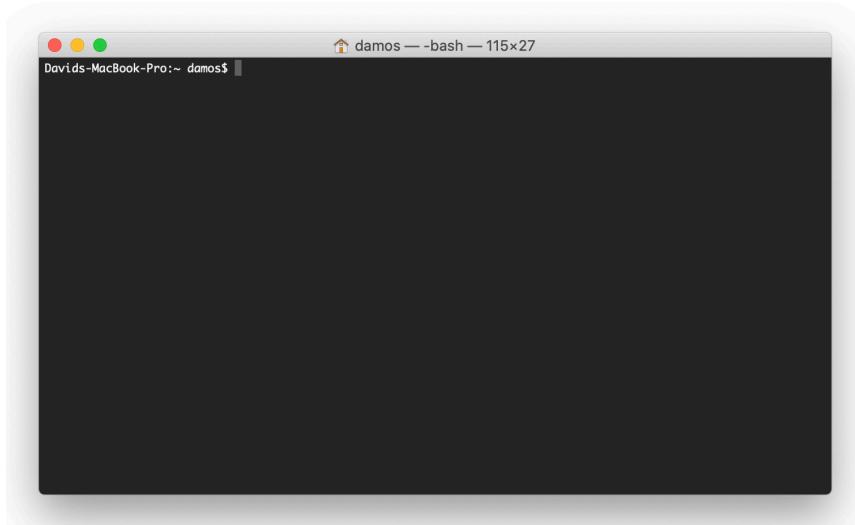


Alternatively, you may use the Powershell application by pressing the Windows key, typing `powershell` and pressing `Enter`. The Powershell window looks like this:



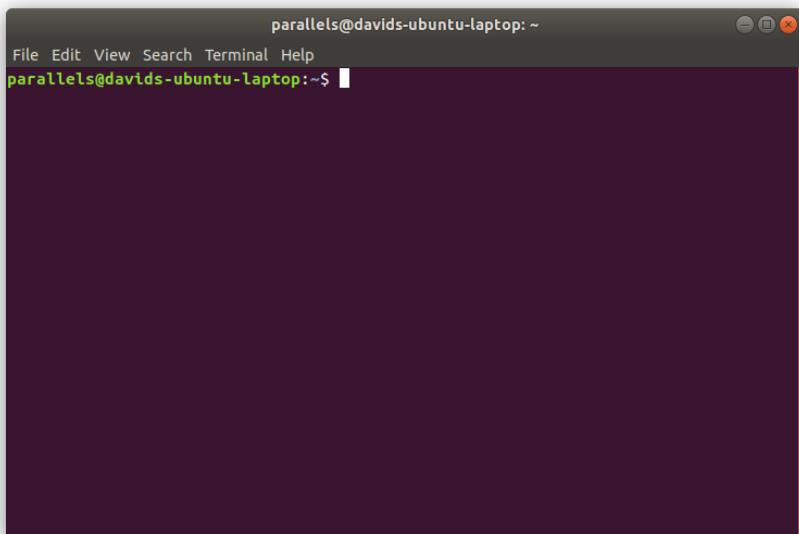
macOS

Press **Cmd** + **Spacebar** to open the Spotlight search window. Type **terminal** and press **Enter** to open the Terminal app. The window that opens look like this:



Ubuntu Linux

Click on the **Show Applications** button at the bottom of your toolbar and search for terminal. Then click on the *Terminal* application icon to open the terminal. The window that opens looks something like this:



With your terminal program open, type in the following command to check whether or not `pip` is installed on your system:

```
$ python3 -m pip --version
```

If `pip` is installed, you should see something like the following output displayed in your terminal:

```
pip 19.3.1 from c:\users\David\appdata\local\programs\python\python38-32\lib\site-packages\pip (python 3.8)
```

This output indicates that version 19.3.1 of `pip` is currently installed and is linked to the Python 3.8 installation.

Important

On macOS and Ubuntu Linux, it is important to run `pip` commands using the `python3` command and not `python`. This ensures that the `pip` installation for Python 3 is used, instead of the Python 2 version of `pip` that may have come pre-installed on your machine.

On Windows, the `python3` command may not work. If you do not see any output, or encounter an error, try the command with `python` instead:

```
$ python -m pip --version
```

If that works for you, replace all the `python3` commands in this chapter with `python`.

Note that the version you see displayed on your computer may be different, and that it might be linked to a different Python installation. This is just fine, as long as the version of Python you see displayed is any version of Python 3.

If your operating system tells you that `pip3` is an unrecognized command, then `pip` was not installed with your Python distribution. In that case, you may want to review the instructions for installing Python in Chapter 2.

Upgrading pip to the Latest Version

Before we go any further, let's make sure that you have the latest version of `pip` installed. To upgrade `pip`, run type the following into your terminal and press `Enter`:

```
$ python3 -m pip install --upgrade pip
```

If a newer version of `pip` is available, it will be downloaded and installed. Otherwise, you will see a message indicating that the latest version is already installed. This message usually says something like:

Requirement already satisfied.

Now that you have `pip` upgraded to the latest version, let's see what we can do with it!

Listing All Installed Packages

You can use `pip` to list all of the packages you have installed. Let's take a peek at what is currently available. Type the following into your terminal:

```
$ python3 -m pip list
```

If you haven't already installed any packages, which should be the case if you started this course with a fresh Python 3.8 installation, you should see something like the following:

Package	Version
-----	-----
pip	19.3.1
setuptools	41.2.0

As you can see, there isn't much here. You see `pip` itself listed, because `pip` is a package. You may also see `setuptools`. This is a package used by `pip` to setup and install other packages.

When you install a package with `pip`, it will show up in this list. You can always use `pip list` to see which packages, and which version of each package, you currently have installed.

Installing a Package

Let's install your first Python package! For this exercise, you will install the `requests` package, which is one of the most popular Python packages ever created. In your terminal, type the following:

```
$ python3 -m pip install requests
```

While `pip` is installing the `requests` package, you will see a bunch of output:

```
Collecting requests
  Downloading https://.../requests-2.22.0-py2.py3-none-any.whl (57kB)
    ..... | 61kB 2.0MB/s
Collecting urllib3!=1.25.0,!=1.25.1,<1.26,>=1.21.1
  Downloading https://.../urllib3-1.25.7-py2.py3-none-any.whl (125kB)
    ..... | 133kB 3.3MB/s
Collecting certifi>=2017.4.17
  Downloading https://.../certifi-2019.11.28.py3-none-any.whl (156kB)
    ..... | 163kB ...
Collecting chardet<3.1.0,>=3.0.2
  Downloading https://.../chardet-3.0.4-py2.py3-none-any.whl (133kB)
    ..... | 143kB 6.8MB/s
Collecting idna<2.9,>=2.5
  Downloading https://.../idna-2.8-py2.py3-none-any.whl (58kB)
    ..... | 61kB 3.8MB/s
Installing collected packages: urllib3, certifi, chardet, idna,
  requests
Successfully installed certifi-2019.11.28 chardet-3.0.4 idna-2.8
  requests-2.22.0 urllib3-1.25.7
```

Note

The formatting of the above output has been altered so that it fits nicely on the page. The output that you see on your computer may look different.

Notice that `pip` first tells you that it is “Collecting `requests`.” You will see the URL that `pip` is using to install the package from, as well as a progress bar indicating the progress of the download.

After that, you will see that `pip` installs four more packages: `chardet`, `certifi`, `idna` and `urllib3`. These packages are **dependencies** of `requests`. That means that `requests` requires these packages to be installed in order for it to work properly.

Once `pip` is done installing `requests` and its dependencies, run `pip3 list` in your terminal again. You should now see the following list:

```
$ python3 -m pip list
Package      Version
-----
certifi      2019.11.28
chardet      3.0.4
idna         2.8
pip          19.3.1
requests     2.22.0
setuptools   41.2.0
urllib3      1.25.7
```

You can see that version 2.22.0 of `requests` was installed, as well as the `chardet`, `certifi`, `idna`, and `urllib3` dependencies.

By default, `pip` installs the latest version of a package. You can control which version of a package gets installed with some optional version specifiers.

Installing Specific Package Versions With Version Specifiers

There are several ways to control which version of a package gets installed. For example, you can:

1. Install the latest version greater than some version number
2. Install the latest version less than some version number
3. Install a specific version number

To install the latest version of `requests` whose version number is 2 or greater, you can execute the following:

```
$ python3 -m pip install requests>=2.0
```

Notice the `>=2.0` after the package name `requests`. This tells `pip` to install the latest version of `requests` that is greater than or equal to ver-

sion 2.0.

The symbol `<=` is called a **version specifier** because it specifies which version of the package should be installed. There are several different version specifiers that you can use. Here are the most widely used ones:

Version Specifier	Description
<code><=, >=</code>	Inclusive less than and greater than specifiers
<code><, ></code>	Exclusive less than and greater than specifiers
<code>==</code>	Exactly equal to specifier

Let's look at some examples.

To install the latest version that is less than or equal to some number, use the `<=` version specifier:

```
$ python3 -m pip install requests<=3.0
```

This will install the latest version of `requests` that is less than or equal to version 3.0.

The `<=` and `>=` version specifiers are **inclusive** because they include the version number that follows the specifier. **Exclusive** versions, `<` and `>`, exist as well.

For instance, the following command installs the latest version of `requests` that is strictly less than version 3.0:

```
$ python3 -m pip install requests<3.0
```

You can combine version specifiers to ensure `pip` installs the latest version within a specified version range. For example, the following command installs the latest version of `requests` in the 1.0 series:

```
$ python3 -m pip install requests>=1.0,<2.0
```

You would use something like the above command if your project was

only compatible with the 1.0 series of the package and you want to make sure you install the latest updates to that series.

Finally, you can **pin** dependencies to a specific version with the == version specifier:

```
$ python3 -m pip install requests==2.22.0
```

This command installs exactly version 2.22.0 of the `requests` package.

Show Package Details

Now that you've installed the `requests` package, you can use `pip` to view some details about the package:

```
$ python3 -m pip show requests
Name: requests
Version: 2.22.0
Summary: Python HTTP for Humans.
Home-page: http://python-requests.org
Author: Kenneth Reitz
Author-email: me@kennethreitz.org
License: Apache 2.0
Location: c:\users\David\...\python\python38-32\lib\site-packages
Requires: chardet, idna, certifi, urllib3
Required-by:
```

The `python3 -m pip show` command displays information about an installed package, including the author's name and email, and a home page you can navigate to in your internet browser to learn more about what the package does.

The `requests` package is used for making HTTP requests from a Python program. It is extremely useful in a variety of domains, and is a dependency of a large number of other Python packages.

Uninstalling a Package

If you can install a package with pip, it only makes sense that you can also uninstall a package. Let's uninstall the requests package now.

To uninstall requests, type the following into your terminal:

```
$ python3 -m pip uninstall requests
```

Important

If you already have projects that use requests or one of its dependencies, you may not want to run the commands in the remainder of this section.

You will immediately see the following prompt:

```
Uninstalling requests-2.22.0:  
Would remove:  
  c:\users\damos\...\requests-2.22.0.dist-info\*  
  c:\users\damos\...\requests\*  
Proceed (y/n)?
```

Before pip actually removes anything from your computer, it asks for your permission first. How considerate!

Type y and press **Enter** to continue. You should then see the following message confirming that requests was removed:

```
Successfully uninstalled requests-2.22.0
```

Take a look at your package list again:

```
$ python3 -m pip list  
Package      Version  
-----  
certifi      2018.4.16  
chardet      3.0.4
```

```
idna      2.7
pip       10.0.1
setuptools 39.0.1
urllib3    1.23
```

Notice that `pip` uninstalled `requests`, but it didn't remove any of its dependencies! This behavior is a feature, not a bug.

Imagine that you have installed several packages into your environment with `pip`, some of which share dependencies. If `pip` uninstalled a package *and* its dependencies, it would render any other package requiring those dependencies unusable!

For now, though, go ahead and remove the remaining packages by running `pip uninstall`. You can uninstall all four packages in a single command:

```
$ python3 -m pip uninstall certifi chardet idna urllib3
```

:::

When you are done, verify that everything has been removed by running `pip list` again. You should see the same list of packages you saw when you first started:

```
Package      Version
-----
pip          10.0.1
setuptools   39.0.1
```

Python's ecosystem of third-party packages is one of its greatest strengths. These packages allow Python programmers to be highly productive and create full-featured software much more quickly than can be done in, say, a language like C++.

That said, using third-party packages in your code introduces several concerns that must be addressed with care. You'll learn about some of the pitfalls associated with third-party packages in the next section.

[Leave feedback on this section »](#)

13.2 The Pitfalls of Third-Party Packages

The beauty of third-party packages is that they give you the ability to add functionality to your project without having to implement everything from scratch. This offers massive boosts in productivity.

But with great power comes great responsibility. As soon as you include someone else's package in your project, you are placing an enormous amount of trust in those responsible for developing and maintaining the package.

By using a package you did not develop, you lose control over certain aspects of your project. In particular, the maintainers of a package may release a new version that introduces changes that are incompatible with the version you use in your project.

By default, `pip` installs the latest release of a package, so if you distribute your code to someone else and they install a newer version of a package required by your project, they may not be able to run your code.

This presents a significant challenge, for both the end user and yourself. Fortunately, Python comes with a fix for this all-to-common problem: virtual environments.

A virtual environment creates an isolated and, most importantly, reproducible environment that you can use to develop a project. The environment can contain a specific version of Python, as well as specific versions of your project's dependencies.

When you distribute your code to someone else, they can reproduce this environment and be confident that they can run your code without error.

Virtual environments are a more advanced topic outside the scope of this book. To learn more about virtual environments and how to use them, check out Real Python’s [Managing Python Dependencies With Pip and Virtual Environments](#) course. In it you will learn how to:

- Install, use, and manage third-party Python packages with the “pip” package manager on Windows, macOS, and Linux, in more detail than presented here.
- Isolate project dependencies with so-called virtual environments to avoid version conflicts in your Python projects.
- Apply a complete 7-step workflow for finding and identifying quality third-party packages to use in your own Python projects (and justifying your decisions to your team or manager.)
- Set up repeatable development environments and application deployments using the “pip” package manager and requirements files.

[Managing Python Dependencies With Pip and Virtual Environments](#) is a great next step when you have completed this book.

[Leave feedback on this section »](#)

13.3 Summary and Additional Resources

In this chapter, you learned how to install third-party packages using Python’s package manager `pip`. You saw several useful `pip` commands, including `pip install`, `pip list`, `pip show` and `pip uninstall`.

You also learned about some of the pitfalls associated with third party packages. Not every package that is downloadable with `pip` is a good choice for your project. Since you do not have control over the code in the package you install, you must trust that the package is safe and will work well for the users of your program.

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-12

Additional Resources

To learn more about managing third-party packages, you can check out these resources:

- [Managing Python Dependencies Course](#)
- [Python Virtual Environments: A Primer](#)
- Recommended resources on realpython.com

[Leave feedback on this section »](#)

Chapter 14

Creating and Modifying PDF Files

The **PDF**, or **portable document format**, is one of the most common formats for sharing documents over the internet. PDF files can contain text, images, tables, forms, and even rich media like videos and animations, all in a single file.

The abundance of content types PDFs may contain can make working with them somewhat difficult. There's a lot of data to decode when opening a PDF file! Fortunately, the Python ecosystem has some great packages for reading, manipulating, and creating PDF files!

In this chapter, you will learn how to:

- Read text from a PDF
- Extract pages and split a PDF into multiple files
- Concatenate and Merge PDF files
- Rotate and crop pages in a PDF file
- Encrypt and Decrypt PDF files with passwords
- Create a PDF file from scratch

Let's get started!

[Leave feedback on this section »](#)

14.1 Extract Text From a PDF

In this section, you'll learn how to read a PDF file and extract the text using the [PyPDF2](#) package. Before you can do that, though, you need to install PyPDF2 with pip:

```
$ python3 -m pip install PyPDF2
```

Verify the installation by running the following in your terminal:

```
$ python3 -m pip show PyPDF2
Name: PyPDF2
Version: 1.26.0
Summary: PDF toolkit
Home-page: http://mstamy2.github.com/PyPDF2
Author: Mathieu Fenniak
Author-email: biziqe@mathieu.fenniak.net
License: UNKNOWN
Location: c:\users\david\python38-32\lib\site-packages
Requires:
Required-by:
```

Pay particular attention to the version information. At the time of writing, the latest version of PyPDF2 is 1.26.0. You'll need to restart IDLE if you have it open in order to use the [PyPDF2](#) package.

Now that you have PyPDF2 installed let's start working with some PDF files!

Open a PDF File

Let's get started by opening a PDF and reading some information about it. We'll use the `Pride_and_Prejudice.pdf` file located in the Chapter 14 Practice Files folder.

Note

If you haven't already, you can download the exercise solutions and practice files [here](#).

Open IDLE's interactive window and import the PdfFileReader class from the PyPDF2 package:

```
>>> from PyPDF2 import PdfFileReader
```

To create a new instance of the PdfFileReader class, you'll need to path to the PDF file that you want to open. Let's get that now using the pathlib module:

```
>>> from pathlib import Path  
>>> pdf_path = (  
...     Path.home() /  
...     "python-basics-exercises" /  
...     "ch13-interact-with-pdf-files" /  
...     "practice_files" /  
...     "Pride_and_Prejudice.pdf"  
... )
```

The pdf_file path variable now contains the path to a PDF version of Jane Austen's *Pride and Prejudice*.

Note

You may need to change pdf_path so that it corresponds to the location of the python-basics-exercises/ folder on your computer.

Now create the PdfFileReader instance:

```
>>> pdf = PdfFileReader(str(pdf_path))
```

pdf_path is converted to a string because PdfFileReader doesn't know how to read from a pathlib.Path object.

Recall from *Chapter 12: File Input and Output* that all open files

should be closed before a program terminates. The `PdfFileReader` object does all of this for you, so you don't need to worry about opening or closing the PDF file!

Now that you've created a `PdfFileReader` instance, you can use it to gather information about the PDF. For instance, the `.getNumPages()` method returns the number of pages contained in the PDF file:

```
>>> pdf.getNumPages()  
234
```

The `Pride_and_Prejudice.pdf` file has 234 pages!

Notice that `.getNumPages()` is written in camelCase and not `snake_case`, as recommended in [PEP 8](#). Remember, PEP 8 is a set of guidelines, not rules. As far as Python is concerned, camelCase is perfectly acceptable.

Note

PyPDF2 is adapted from the PyPDF package. PyPDF was written in 2005, only four years after PEP 8 was published.

At that time, many Python programmers were migrating from languages where camelCase is more common.

You can also access some document information using the `.documentInfo` attribute:

```
>>> pdf.documentInfo  
{'/Title': 'Pride and Prejudice, by Jane Austen', '/Author': 'Chuck', '/Creator':
```

The object returned by `.documentInfo` looks like a dictionary, but it's not really the same thing. You can access each item in `.documentInfo` as an attribute.

For example, to get the title, use the `.title` attribute:

```
>>> pdf.documentInfo.title  
'Pride and Prejudice, by Jane Austen'
```

The `.documentInfo` object contains the PDF **metadata** which is set when a PDF is created.

The `PdfFileReader` class is the gateway to working with PDF files in Python. It provides all the necessary methods and attributes needed to access data in a PDF file.

Let's explore what you can do with a PDF file and how you do it!

Extract Text From a Page

PDF pages are represented in PyPDF2 with the `PageObject` class. You use `PageObject` instances to interact with pages in a PDF file.

You don't need to create your own `PageObject` instances. Instead, you access them through a `PdfFileReader` object. Let's see how this is done by extracting the text from the first page of the `Pride_and_Prejudice.pdf` file.

There are two steps to extracting text from a single PDF page:

1. Get a `PageObject` with `PdfFileReader.getPage()`
2. Extract the text as a string with the `PageObject` instance's `.extractText()` method.

`Pride_and_Prejudice.pdf` has 243 pages. Each page has an index between 0 and 242. You can get an object representing a specific page by passing the page's index to the `PdfFileReader.getPage()` method:

```
>>> first_page = pdf.getPage(0)
```

`.getPage()` returns a `PageObject`:

```
>>> type(first_page)  
<class 'PyPDF2.pdf.PageObject'>
```

You can extract the page's text with the `PageObject.extractText()` method:

```
>>> first_page.extractText()
'\n \nThe Project Gutenberg EBook of Pride and Prejudice, by Jane
Austen\n \n\nThis eBook is for the use of anyone anywhere at no cost
and with\n \nalmost no restrictions whatsoever. You may copy it,
give it away or\n \nre\n \nuse it under the terms of the Project
Gutenberg License included\n \nwith this eBook or online at
www.gutenberg.org\n \n \nTitle: Pride and Prejudice\n \n
\nAuthor: Jane Austen\n \n \nRelease Date: August 26, 2008
[EBook #1342]\n\n[Last updated: August 11, 2011]\n \n \nLanguage:
Eng\nlish\n \n \nCharacter set encoding: ASCII\n \n \n***\n
START OF THIS PROJECT GUTENBERG EBOOK PRIDE AND PREJUDICE ***\n \n
\n \n \n \nProduced by Anonymous Volunteers, and David Widger\n
\n \n \n \n \n \nPRIDE AND PREJUDICE \n \n \nBy Jane
Austen \n \n\n \n \nContents\n \n '
```

Note that the output displayed here has been formatted to fit better on this page. The output you see on your computer may be formatted differently.

To extract all of the text from the entire PDF, you'll need to away to iterate over all of the pages in the document.

Every `PdfFileReader` object has a `.pages` attribute used to access an iterable of `PageObject` objects for each page in the PDF. This iterable is in order, so the first `PageObject` corresponds to the first page of the PDF, the second `PageObject` to the second page, and so on.

Here's how you use a `for` loop to loop over all the pages in the PDF and print their text:

```
>>> for page in pdf.pages:
...     print(page.extractText())
...
```

Let's combine everything you've learned and write a program that ex-

tracts all of the text from the `Pride_and_Prejudice.pdf` file and saves it to a `.txt` file.

Putting It All Together

Open a new script window in IDLE. Type out the script below:

```
from pathlib import Path
from PyPDF2 import PdfFileReader

# Change the path below to the correct path for your computer.
pdf_path = (
    Path.home() /
    "python-basics-exercises" /
    "ch13-interact-with-pdf-files" /
    "practice-files" /
    "Pride_and_Prejudice.pdf"
)

# 1
pdf_reader = PdfFileReader(str(pdf_path))
output_file_path = Path.home() / "Pride_and_Prejudice.txt"

# 2
with output_file_path.open(mode="w") as output_file:
    # 3
    output_file.write(
        f"{pdf_reader.documentInfo.title}\n"
        f"Number of pages: {pdf_reader.getNumPages()}\n\n"
    )

# 4
for page in pdf_reader.pages:
    text = page.extractText()
    output_file.write(text)
```

Let's break that down.

First, you assign a new `PdfFileReader` instance to the `pdf_reader` variable and a new `Path` object to the file `Pride_and_Prejudice.txt` in your home directory to the `output_file_path` variable (#1).

Next, the script opens the `output_file_path` in write mode (#3). The file object returned by `.open()` is assigned to the variable `output_file`.

The `with` statement, which you learned about in *Chapter 12: File Input and Output*, ensures that the file is properly closed when the code in indented `with` block finished executing.

Inside the `with` block, the PDF title and number of pages are written to the text file using `output_file.write()` (#3). After that, each `PageObject` in the `pdf_reader.pages` iterable is looped over in a `for` loop (#4).

At each step in the `for` loop, the `page` variable is assigned to the next `PageObject` in the iterable. Then the text from each page is extracted with `page.extractText()` and written to the `output_file`.

When you save and run the above script, a new file called `Pride_and_Prejudice.txt` containing the full text of the `Pride_and_Prejudice.pdf` document is created in your home directory. Open it up and check it out!

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. In the Chapter 14 Practice Files directory there is a PDF file called `zen.pdf`. Create a `PdfFileReader` from this PDF.
2. Using the `PdfFileReader` instance from Exercise 1, print the total number of pages in the PDF.
3. Print the text from the first page of the PDF in Exercise 1.

[Leave feedback on this section »](#)

14.2 Extract Pages From a PDF

In the last section, you learned how to extract all of the text from a PDF file and save the text to a .txt file. Now you'll learn how to extract a page, or a range of pages, from an existing PDF and save them to a new PDF.

The PdfFileWriter class is used to created a new PDF file. Let's explore this class and learn the steps needed to create a PDF using PyPDF2.

The PdfFileWriter Class

The PdfFileWriter class is used to create new PDF files. In IDLE's interactive window, import the PdfFileWriter class and create a new instance called pdf_writer:

```
>>> from PyPDF2 import PdfFileWriter  
>>> pdf_writer = PdfFileWriter()
```

PdfFileWriter objects are containers for pages. To create an new PDF, you need to add some PageObject instances to the PdfFileWriter and then write those pages to a file.

Let's add a blank page to the pdf_writer object:

```
>>> page = pdf_writer.addBlankPage(width=72, height=72)
```

The .addBlankPage() method adds a blank page to the PDF writer object. Since there are no pages in the writer, it is added as the first page.

The width and height parameters are required and determine the width and height of the page in **points**. One point is equal to 1/72 inches. So the above code adds a one inch square blank page to pdf_writer.

.addBlankPage() returns a new PageObject instance representing the page that was added to the PdfFileWriter:

```
>>> type(page)
<class 'PyPDF2.pdf.PageObject'>
>>> page
{'/Type': '/Page', '/Parent': IndirectObject(1, 0), '/Resources': {}, '/MediaBox': RectangleObject([0, 0, 72, 72])}
```

In this example you've assigned the `PageObject` instance returned by `.addBlankPage()` to the `page` variable, but in practice you don't usually need to do this. That is, you usually call `.addBlankPage()` without assigning the return value to anything:

```
>>> pdf_writer.addBlankPage(width=72, height=72)
```

With at least one page added to `pdf_writer`, you can write the contents to a new PDF file. To do this, pass a file object in binary write mode to the `PdfFileWriter.write()` method:

```
>>> from pathlib import Path
>>> with Path("blank.pdf").open(mode="wb") as output_file:
...     pdf_writer.write(output_file)
...
>>>
```

This creates a new file called `blank.pdf` in your current working directory. If you open the file with a PDF reader, such as Adobe Acrobat, you'll see a document with a single blank one inch square page.

Important

Notice that you save the PDF file by passing the file object to the `PdfFileWriter` object's `.write()` method, and *not* the file object's `.write()` method.

In particular, the following code will not work:

```
>>> with Path("blank.pdf").open(mode="wb" as output_file):
...     output_file.write(pdf_writer)
```

This might seem backwards to you, so make sure you avoid this mistake!

`PdfFileWriter` objects can write to new PDF files, but they can't create new content from scratch other than blank pages. This might seem like a big problem, but in many situations you don't need to create new content. Often, you'll work with pages extracted from PDF files that you've opened with a `PdfFileReader` instance.

Note

You'll learn how to create PDF files from scratch in Section 13.8 *Create a PDF File From Scratch*.

In the example you saw above, there were three steps to create a new PDF file using PyPDF2:

1. Create a `PdfFileWriter` instance
2. Add one or more pages to the `PdfFileWriter` instance
3. Write to a file using the `PdfFileWriter.write()` method

You'll see this pattern over and over as you learn various ways to add pages to a `PdfFileWriter` instance.

Extracting a Single Page From a PDF

Let's revisit the *Pride and Prejudice* PDF that you worked with in the last section. We'll open the PDF using a `PdfFileReader` class instance, extract the first page of the PDF, and then create a new PDF file containing just the single extracted page.

Open IDLE's interactive window and import both `PdfFileReader` and `PdfFileWriter` from `PyPDF2`, as well as the `Path` class from the `pathlib` module:

```
>>> from pathlib import Path  
>>> from PyPDF2 import PdfFileReader, PdfFileWriter
```

Now open the `Pride_and_Prejudice.pdf` file with a `PdfFileReader` instance:

```
>>> pdf_path = (  
...     Path.home() /  
...     "python-basics-exercises" /  
...     "ch13-interact-with-pdf-files" /  
...     "practice_files" /  
...     "Pride_and_Prejudice.pdf"  
... )  
>>> input_pdf = PdfFileReader(str(pdf_path))
```

Remember, you may need to change the file path so that it works on your system.

The first page in the PDF is at index 0. You can extract it as a `PageObject` by passing the argument 0 to `.getPage()`

```
>>> first_page = input_pdf.getPage(0)
```

Now you can create a new `PdfFileWriter` instance and add the extracted page to it:

```
>>> pdf_writer = PdfFileWriter()  
>>> pdf_writer.addPage(first_page)
```

The `.addPage()` method adds a page to the set of pages in the `pdf_writer` object, just like `.addBlankPage()`. The difference is that you must pass a `PageObject` to `.addPage()`.

Now write the contents of `pdf_writer` to a new file called `first_page.pdf`:

```
>>> with Path("first_page.pdf").open(mode="wb") as output_file:  
...     pdf_writer.write(output_file)  
...  
>>>
```

You now have a new PDF file saved in your current working directory with the name `first_page.pdf` that contains the cover page of the `Pride_and_Prejudice.pdf` file. Pretty neat!

Extract Multiple Pages From a PDF File

Using `for` loops, you can extract multiple pages from a PDF file and save them to a new PDF. Let's extract the first chapter from `Pride_and_Prejudice.pdf`.

If you open `Pride_and_Prejudice.pdf` with a PDF viewer, you can see that the first chapter is on the second, third, and fourth pages in the PDF. Since pages are indexed starting with 0, we'll need to extract the pages at the indices 1, 2, and 3.

Let's set everything up by importing the classes we need and opening the PDF file:

```
>>> from PyPDF2 import PdfFileReader, PdfFileWriter  
>>> from pathlib import Path  
>>> pdf_path = (  
...     Path.home() /  
...     "python-basics-exercises" /  
...     "ch13-interact-with-pdf-files" /  
...     "practice_files" /  
...     "Pride_and_Prejudice.pdf"  
... )
```

```
>>> input_pdf = PdfFileReader(str(pdf_path))
```

Our goal is to extract the pages at indices 1, 2, and 3, add these to a new `PdfFileWriter` instance, and then write them to a new PDF file.

One way to do this is to loop over the range of numbers starting at 1 and ending at 3, extracting the page at each step of the loop and adding it to the `PdfFileWriter` instance. Here's what that looks like in code:

```
>>> pdf_writer = PdfFileWriter()
>>> for n in range(1, 4):
...     page = input_pdf.getPage(n)
...     pdf_writer.addPage(page)
...
>>>
```

Now `pdf_writer` has three pages added to it, which you can check with the `.getNumPages()` method:

```
>>> pdf_writer.getNumPages()
3
```

Finally, you can write the extracted pages to a new PDF file:

```
>>> with Path("chapter1.pdf").open(mode="wb") as output_file:
...     pdf_writer.write(output_file)
...
>>>
```

Now you can open the `chapter1.pdf` file in your current working directory to read just the first chapter of *Pride and Prejudice*.

Another way to extract multiple pages from a PDF is to take advantage of the fact that the iterable returned by `PdfFileReader.pages` supports slice notation.

Let's redo the previous example using `.pages` instead of looping over a range object. Since you've already imported the necessary classes and set-up to file paths, start by initializing a new `PdfFileWriter` object:

```
>>> pdf_writer = PdfFileWriter()
```

Now loop over a slice of the `.pages` iterable from indices starting at 1 and ending at 3:

```
>>> for page in input_pdf.pages[1:4]:  
...     pdf_writer.addPage(page)  
...  
>>>
```

Recall that the values in a slice range from the item at the first index in the slice and up to, but not including, the item at the second index in the slice. So `.pages[1:4]` returns an iterable of the pages with indices 1, 2, and 3.

Finally, write the contents of `pdf_writer` to the output file:

```
>>> with Path("chapter1_slice.pdf").open(mode="wb") as output_file:  
...     pdf_writer.write(output_file)  
...  
>>>
```

Now open the `chapter1_slice.pdf` file in your current working directory and compare it to the `chapter1.pdf` file you made by looping over the `range` object. They contain the same pages!

Sometimes you need to extract every page from a PDF. You can use the methods illustrated above to do this, but PyPDF2 provides a shortcut. `PdfFileWriter` instances have an `.appendPagesFromReader()` method that is used to append pages from a `PdfFileReader` instance.

To use `.appendPagesFromReader()` pass a `PdfFileReader` instance to its `reader` parameter. For example, the following copies every page from the *Pride and Prejudice* PDF to a `PdfFileWriter` instance:

```
>>> # Assume pdf_reader contains the opened Pride_and_Prejudice.pdf  
>>> pdf_writer = PdfFileWriter()  
>>> pdf_writer.appendPagesFromReader(pdf_reader)
```

pdf_writer now contains every page in pdf_reader!

Review Exercises

1. Extract the last page from the `Pride_and_Prejudice.pdf` file and save it to a new file called `last_page.pdf` in your home directory.
2. Extract all pages with even numbered *indices* from the `Pride_and_Prejudice.pdf` and save them to a new file called `every_other_page.pdf` in your home directory.
3. Split the `Pride_and_Prejudice.pdf` file into two new PDF files. The first file should contain the first 150 pages, and the second file should contain the remaining pages. Save both files in your home directory as `part_1.pdf` and `part_2.pdf`.

[Leave feedback on this section »](#)

14.3 Challenge: PdfFileSplitter Class

Create a class called `PdfFileSplitter` that reads a PDF from an existing `PdfFileReader` instance and splits the PDF into two new PDFs.

The class should be instantiated with a path string. For example, here's how you would create a `PdfFileSplitter` instance from a PDF called `mydoc.pdf` in your current working directory:

```
pdf_splitter = PdfFileSplitter("mydoc.pdf")
```

The `PdfFileSplitter` class should have two methods:

1. `.split()` that has a single parameter `breakpoint` that expects an integer representing the page number to split the PDF.

After `.split()` is called, the `PdfFileSplitter` class should have an attribute `.writer1` assigned to a `PdfFileWriter` instance containing all the pages in the original PDF up to *but not including* the breakpoint page, and `.writer2` assigned to a `PdfFileWriter` instance containing the remaining pages in the original PDF.

- .write() that has a single parameter filename that expects a path string.

When .write() is called, two PDFs should be written to the specified path. The first one with the name filename + "_1.pdf" and the second with the name filename + "_2.pdf".

For example, here's how you would split the mydoc.pdf at page four:

```
pdf_splitter.split(breakpoint=4)
```

Then, to write two new PDFs in the current working directory as mydoc_split_1.pdf and mydoc_split_2.pdf, you would call .write() with the file name "mydoc_split":

```
pdf_splitter.write("mydoc_split")
```

Check that the splitter works by splitting the Pride_and_Prejudice.pdf file in the Chapter 14 Practice Files folder with the breakpoint at the 150th page.

[Leave feedback on this section »](#)

14.4 Concatenating and Merging PDFs

Two common tasks when working with PDF files are concatenating and merging several PDFs together into a single file.

When you concatenate two or more PDFs together, you join the files into a single document one after another. For example, a company may concatenate several daily reports into one monthly report at the end of a month.

Merging two PDFs together also joins two PDFs into a single file, but instead of joining the second PDF at the end of the first, it can be inserted after a specific page in the first PDF, pushing all following pages in the first PDF to the end of the second one.

In this section, you'll learn how to concatenate and merge PDFs using the `PyPDF2` package's `PdfFileMerger`.

The `PdfFileMerger` Class

The `PdfFileMerger` class is a lot like the `PdfFileWriter` class you learned about in the last section. Both classes are used to write PDF files. In both cases, pages are added to instances of the class and then written to a file.

The main differences is that `PdfFileWriter` can only append pages to the end of the list of pages already contained in the writer, whereas `PdfFileMerge` can insert pages at any location.

Let's go ahead and create our first `PdfFileMerger` instance. In IDLE's interactive window, type the following:

```
>>> from PyPDF2 import PdfFileMerger  
>>> pdf_merger = PdfFileMerger()
```

First, import the `PdfFileMerger` class from the `PyPDF2` package. Then create a new `PdfFileMerger` instance and assigns it to the `pdf_merger` variable. `PdfFileMerger` objects are empty when they are first instantiated. We'll need to add some pages to it before we can do anything.

There are a couple of ways to add pages to the `pdf_merger` object, and how you do that depends on what you need to accomplish:

- `.append()` concatenates every page in an existing PDF document to the end of the pages currently in the `PdfFileMerger`.
- `.merge()` is used to insert all of the pages in an existing PDF document after a specific page in the `PdfFileMerger`.

We'll look at both methods in this section, starting with `.append()`.

Concatenating PDFs With `.append()`

In the Chapter 14 Practice Files directory of the [Python Basics Exercises](#) repository, there is a subdirectory called `expense_reports` with three expense reports for an employee named Peter Python.

Peter needs to concatenate these three PDFs together and submit them to his employer as a single PDF file so that he can get reimbursed for some work-related expenses.

Let's start by using the `pathlib` module to get a list `Path` objects for each of the three expense reports in the `expense_reports/` folder:

```
>>> from pathlib import Path  
>>> reports_dir = (  
...     Path.home() /  
...     "python-basics-exercises" /  
...     "ch13-interact-with-pdf-files" /  
...     "practice_files" /  
...     "expense_reports"  
... )
```

After you import the `Path` class, you need to build the path to the `expense_reports/` directory. Note that you may need to alter the code above to get the correct path on your computer.

Once you have the path to the `expense_reports/` directory assigned to the `reports_dir` variable, you can use `.glob()` to get an iterable of paths to PDF files in the directory.

Let's take a look at what's in the directory:

```
>>> for path in reports_dir.glob("*.pdf"):  
...     print(path.name)  
...  
Expense report 1.pdf  
Expense report 3.pdf  
Expense report 2.pdf
```

The names of the three files are listed, but they aren't in order. Furthermore, the order of the files you see in the output on your computer may not match the output shown here.

In general, the order of paths returned by `.glob()` is not guaranteed, so you'll need to order them yourself. You can do this by creating a list containing the three file paths, and then calling the `.sort()` method on that list:

```
>>> expense_reports = list(reports_dir.glob("*.pdf"))
>>> expense_reports.sort()
```

Recall that `.sort()` sorts a list in place, so you don't need to assign the return value to a variable. After calling `.sort()`, the `expense_reports` list is sorted by file name in alphabetical order.

Let's check that the sorting worked by looping over `expense_reports` again and printing out the file names:

```
>>> for path in expense_reports:
...     print(path.name)
...
Expense report 1.pdf
Expense report 2.pdf
Expense report 3.pdf
```

That looks good!

Now we can concatenate the three PDFs together. To do that, we'll use the `PdfFileMerger.append()` method, which requires a single string argument representing the path to a PDF file. When you call `.append()`, all of the pages in the PDF file are appended to the set of pages in the `PdfFileMerger` object.

Let's see this in action. First import the `PdfFileMerger` class and create a new instance:

```
>>> from PyPDF2 import PdfFileMerger  
>>> pdf_merger = PdfFileMerger()
```

Now loop over the paths in the sorted `expense_reports` list and append them to `pdf_merger`:

```
>>> for path in expense_reports:  
...     pdf_merger.append(str(path))  
...  
>>>
```

Notice that each `Path` object in `expense_reports` is converted to a string with `str()` before being passed to `pdf_merger.append()`.

With all of the PDF files in the `expense_reports/` directory concatenated together in the `pdf_merger` object, the last thing you need to do is write everything to an output PDF file. `PdfFileMerger` instances have a `.write()` method that works just like the `PdfFileWriter.write()`.

Open a new file in binary write mode, then pass the file object to `pdf_merger.write()`:

```
>>> with Path("expense_reports.pdf").open(mode="wb") as output_file:  
...     pdf_merger.write(output_file)  
...  
>>>
```

You now have a PDF file called `expense_reports.pdf` in your current working directory. Open it up with a PDF reader and you'll find all three expense reports together in the same PDF file.

Merging PDFs With `.merge()`

To merge two or more PDFs together, use the `PdfFileMerger.merge()` method. This method is similar to the `.append()` method, except that you must specify where in the output PDF to insert all of the content of the PDF you are merging.

Let's look at an example. Goggle, Inc has prepared a quarterly report,

but forgot to include a table of contents. Peter Python noticed the mistake and quickly created a PDF with the missing table of contents. Now he needs to merge that PDF into the original report.

The first thing you need to do is import everything you need from PyPDF2 and pathlib:

```
>>> from pathlib import Path  
>>> from PyPDF2 import PdfFileMerger
```

The report PDF and table of contents PDF can be found in the `quarterly_report/` subfolder of the Chapter 14 Practice Files folder. The report is in a file called `report.pdf`, and the table of contents is in a file a `toc.pdf`.

Let's go ahead and get the paths to both of those files:

```
>>> report_dir = (  
...     Path.home() /  
...     "python-basics-exercises" /  
...     "ch13-interact-with-pdf-files" /  
...     "practice_files" /  
...     "quarterly_report"  
... )  
>>> report_path = report_dir / "report.pdf"  
>>> toc_path = report_dir / "toc.pdf"
```

The first thing we'll do is append the report PDF to a new `PdfFileMerger` instance using the `.append()` method:

```
>>> pdf_merger = PdfFileMerger()  
>>> pdf_merger.append(str(report_path))
```

Now that `pdf_merger` has some pages in it, we can merge the table of contents PDF into it at the correct location. If you open the `report.pdf` file with a PDF reader, you'll see that the first page of the report is a title page. The second is an introduction, and the remaining pages have different report sections in them.

We want to insert the table of contents after the title page and just before the introduction section. Since PDF pages are indexed starting with 0 in PyPDF2, this means that we need to insert the table of contents after the page at index 0 and before the page at index 1.

To do that, call `pdf_merger.merge()` with two arguments: first the integer 1 indicating the index of the page where the table of contents should be inserted, and second a string containing the path the PDF file for the table of contents.

Here's what that looks like:

```
>>> pdf.merge(1, str(toc_path))
```

Every page in the table of contents PDF is inserted *before* the page at index 1. Since the table of contents PDF is only one page, it gets inserted at index 1 and the page currently at index 1 gets shifted to index 2, the page currently at index 2 gets shifted to index 3, and so on.

Now write the merged PDF to an output file:

```
>>> with Path("full_report.pdf").open(mode="wb") as output_file:  
...     pdf_merger.write(output_file)  
...  
>>>
```

You now have a `full_report.pdf` file in your current working directory. Open it up with a PDF reader and check that the table of contents was inserted at the correct spot.

Concatenating and merging PDFs are common operations. While the examples in this section are admittedly somewhat contrived, you can imagine how useful a program would be for merging thousands of PDFs, or for automating routine tasks that would otherwise take a human lots of time to complete.

Review Exercises

1. In the Chapter 14 Practice Files directory there are three PDFs called `merge1.pdf`, `merge2.pdf`, and `merge3.pdf`. Using a `PdfFileMerger` instance, concatenate the two files `merge1.pdf` and `merge2.pdf` using `.append()`.

Save the concatenated PDFs to a file called `concatenated.pdf` in your home directory.

2. With a new `PdfFileMerger` instance, use `.merge()` to merge the file `merge3.pdf` in-between the two pages in the `concatenated.pdf` file you made in exercise 1. Save the new file to your home directory as `merged.pdf`.

The final result should be a PDF with three pages. The first page should have the number 1 on it, the second should have 2, and the third should have 3.

[Leave feedback on this section »](#)

14.5 Rotating and Cropping PDF Pages

So far you've learned how to extract text from PDF files, extract pages, and concatenate and merge PDF files. These are all common operations with PDF files, but PyPDF2 has many other useful features.

In this section, you'll learn how to rotate and crop pages in a PDF file.

Rotating Pages

Let's start by learning how to rotate pages. For this example, we'll use the `ugly.pdf` file in the Chapter 14 Practice Files folder. The `ugly.pdf` file contains a lovely version of Hans Christian Andersen's *The Ugly Duckling*, except that every odd-numbered page is rotated counter-clockwise by ninety degrees.

Let's fix that. In a new IDLE interactive window, start by importing

the PdfFileReader and PdfFileWriter classes from PyPDF2, as well as the Path class from the pathlib module:

```
>>> from pathlib import Path  
>>> from PyPDF2 import PdfFileReader, PdfFileWriter
```

Now create a Path object for the ugly.pdf file:

```
>>> pdf_path = (  
...     Path.home() /  
...     "python-basics-exercises" /  
...     "ch13-interact-with-pdf-files" /  
...     "practice_files" /  
...     "ugly.pdf"  
... )
```

Finally, let's create new PdfFileReader and PdfFileWriter instances:

```
>>> pdf_reader = PdfFileReader(str(pdf_path))  
>>> pdf_writer = PdfFileWriter()
```

Our goal is to create a new PDF file using pdf_writer that has all of the pages in the PDF rotated correctly. The even numbered pages in the PDF are already properly oriented, but the odd numbered pages in the PDF file are rotated counterclockwise by ninety degrees.

To correct the problem, you'll use the PageObject.rotateClockwise() method. This method takes an integer argument, in degrees, and rotates a page clockwise by that many degrees. For example, .rotateClockwise(90) rotates a PDF page clockwise by ninety degrees.

Note

In addition to the .rotateClockwise() method, the PageObject class also has .rotateCounterClockwise() method for rotating pages counterclockwise.

There are several ways you can go about rotating the pages in the PDF.

We'll discuss two different ways of doing it. Both of them rely on the `.rotateClockwise()` method, but they take different approaches determining which pages get rotated.

The first method is to loop over the indices of the pages in the PDF. During each iteration, check if the index corresponds to a page that needs to be rotated and call `.rotateClockwise()` to rotate the page if needed. Then add the page `pdf_writer`.

Here's what that looks like:

```
>>> for n in range(pdf_reader.getNumPages()):
...     page = pdf_reader.getPage(n)
...     if n % 2 == 0:
...         page.rotateClockwise(90)
...     pdf_writer.addPage(page)
...
>>>
```

Notice that the page gets rotated if the index is even. That might seem odd because it's the odd pages in the PDF that are rotated incorrectly. However, page numbers in the PDF start with 1, while page indices start with 0. That means odd PDF pages have even indices.

If that makes your head spin, don't worry! After years of dealing with stuff like this, even professional programmers get tripped up by these sorts of things!

Note

When you execute the `for` loop above, you'll see a bunch of output in IDLE's interactive window. That's because `.rotateClockwise()` returns a `PageObject` instance.

You can ignore this output for the time being. When you execute programs from IDLE's script window, this output won't be visible.

Now that you've rotated all the pages in PDF, you can write the content

of `pdf_writer` to a new file and check that everything worked:

```
>>> with Path("ugly_rotated.pdf").open(mode="wb") as output_file:  
...     pdf_writer.write(output_file)  
...  
>>>
```

You should now have a file called `ugly_rotated.pdf` in your current working directory with the pages from the `ugly.pdf` file all rotated correctly.

The problem with the approach we just used to rotate the pages in the `ugly.pdf` file is that it depends on knowing ahead of time which pages need to be rotated. In a real-world scenario, it isn't practical to go through an entire PDF taking note of the page numbers of pages to rotate.

In fact, you can determine which pages need to be rotated without prior knowledge. Well, *sometimes* you can.

Let's see how by getting a new `PdfFileReader` instance:

```
>>> pdf_reader = PdfFileReader(str(pdf_path))
```

We need to do this because we altered the pages in the old `PdfFileReader` by rotating them. So by creating a new instance, we're starting fresh.

`PageObject` instances maintain a dictionary of values containing information about the page.

```
>>> pdf_reader.getPage(0)  
{'/Contents': [IndirectObject(11, 0), IndirectObject(12, 0),  
IndirectObject(13, 0), IndirectObject(14, 0), IndirectObject(15, 0),  
IndirectObject(16, 0), IndirectObject(17, 0), IndirectObject(18, 0)],  
'/Rotate': -90, '/Resources': {'/ColorSpace': {'/CS1':  
IndirectObject(19, 0), '/CS0': IndirectObject(19, 0)}}, '/XObject':  
{'/Im0': IndirectObject(21, 0)}, '/Font': {'/TT1':
```

```
IndirectObject(23, 0), '/TTO': IndirectObject(25, 0)}, '/ExtGState':  
{'/GS0': IndirectObject(27, 0)}}, '/CropBox': [0, 0, 612, 792],  
'/Parent': IndirectObject(1, 0), '/MediaBox': [0, 0, 612, 792],  
'/Type': '/Page', '/StructParents': 0}
```

Yikes! Mixed in with all that nonsensical looking stuff is a key called `/Rotate`, which you can see on the fourth line of output above. The value of this key is -90.

You can access the `/Rotate` key on a `PageObject` using subscript notation, just like you can on a Python `dict` object:

```
>>> page = pdf_reader.getPage(0)  
>>> page["/Rotate"]  
-90
```

If you look at the `/Rotate` key for the second page in `pdf_reader`, you'll see that it has a value of 0:

```
>>> page = pdf_reader.getPage(1)  
>>> page["/Rotate"]  
0
```

What all this means is that the page at index 0 has a rotation value of -90 degrees, meaning it has been rotated ninety degree counterclockwise. The page at index 1 has a rotation value of 0, so it has not been rotated at all.

If you rotate the first page using `.rotateClockwise()`, the value of `/Rotate` changes from -90 to 0:

```
>>> page = pdf_reader.getPage(0)  
>>> page["/Rotate"]  
-90  
>>> page.rotateClockwise(90)  
>>> page["/Rotate"]  
0
```

Now that we know how to inspect the `/Rotate` key, let's use it to rotate

the pages in the `ugly.pdf` file.

The first thing we need to do is re-initialize our `pdf_reader` and `pdf_writer` objects so that we get a fresh start:

```
>>> pdf_reader = PdfFileReader(str(pdf_path))
>>> pdf_writer = PdfFileWriter()
```

Now write a loop that loops over the pages in the `pdf_reader.pages` iterable, checks the value of `/Rotate`, and rotates the page if that value is `-90`:

```
>>> for page in pdf_reader.pages:
...     if page["/Rotate"] == -90:
...         page.rotateClockwise(90)
...     pdf_writer.addPage(page)
...
>>>
```

Not only is this loop slightly shorter than the loop in the first solution, but it doesn't rely on any prior knowledge of which pages need to be rotated. You could use a loop like this to rotate pages in any PDF without every having to open it up and look at it.

To finish out the solution, write the contents of `pdf_writer` to a new file:

```
>>> with Path("ugly_rotated2.pdf").open(mode="wb") as output_file:
...     pdf_writer.write(output_file)
...
>>>
```

Now you can open the `ugly_rotated2.pdf` in your current working directory and compare it to the `ugly_rotated.pdf` you generated earlier. They should look identical.

Important

One word of warning about the `/Rotate` key: it is not guaranteed to exist on a page.

If the `/Rotate` key doesn't exist, this usually means that the page has not been rotated. However, that isn't always a safe assumption.

If a `PageObject` has no `/Rotate` key, then a `KeyError` is raised when you try to access it. You can catch this exception with a `try...except` block.

The value of `/rotate` may not always be what you expect. For example, if you scan a paper document with the paper page rotated ninety degrees counter clockwise, then the contents of the PDF will appear rotated. However, the `/Rotate` key may have the value 0.

This is one of many things that can make working with PDF files frustrating. Sometimes, you will just need to open a PDF in a PDF reader program and manually figure some things out.

Cropping Pages

Another common operation with PDFs is cropping pages. You might need to do this to split a single page into multiple pages, or to extract just a small portion of a page, such as a signature or a figure.

For example, there is a file called `half_and_half.pdf` located in the `practice_files/` subdirectory of the `ch13-interact-with-pdf-files/` directory. This PDF contains a portion of Hans Christian Andersen's *The Little Mermaid*.

Each page in this PDF has two columns. Let's split each page of this PDF into two pages, one for each column.

To get started, import the `PdfFileReader` and `PdfFileWriter` classes from `PyPDF2`, and the `Path` class from the `pathlib` module:

```
>>> from pathlib import Path  
>>> from PyPDF2 import PdfFileReader, PdfFileWriter
```

Now create a Path object for the half_and_half.pdf file:

```
>>> pdf_path = (  
...     Path.home() /  
...     "python-basics-exercises" /  
...     "ch13-interact-with-pdf-files" /  
...     "practice_files" /  
...     "half_and_half.pdf"  
... )
```

Next, create a new PdfFileReader object and get the first page of the PDF:

```
>>> pdf_reader = PdfFileReader(str(pdf_path))  
>>> first_page = pdf_reader.getPage(0)
```

In order to crop the page, you need to know a little bit more about how pages are structured. PageObject instances like first_page have an attribute .mediaBox that represents a rectangular area that defines the boundaries of a page.

Let's use IDLE's interactive window to explore the .mediaBox before using it to crop the page:

```
>>> first_page.mediaBox  
RectangleObject([0, 0, 792, 612])
```

Notice that the .mediaBox attribute returns a RectangleObject. This is an object defined in the PyPDF2 package and represents a rectangular area on the page.

You'll notice the list [0, 0, 792, 612] of four numbers in the output. The first two numbers are the x- and y-coordinates of the lower left corner of the rectangle. The third number is the width of the rectangle, and the fourth number represents the height of the rectangle.

Note

The width and height of a `RectangleObject` are defined in **points**. One point is equal to 1/72 inches.

So `RectangleObject([0, 0, 792, 612])` represents a rectangular region with the lower left corner at the origin with a height of 792 points, or 11 inches, and a height of 612 points, or 8.5 inches. Those are the dimensions of a standard letter sized page.

A `RectangleObject` has four attributes that return the coordinates of the rectangle's corners: `.lowerLeft`, `.lowerRight`, `.upperLeft`, and `.upperRight`. Just like the width and height values, coordinates are also given in points.

You can use these four properties to get the coordinates of each corner of the `RectangleObject`:

```
>>> first_page.mediaBox.lowerLeft  
(0, 0)  
>>> first_page.mediaBox.lowerRight  
(792, 0)  
>>> first_page.mediaBox.upperLeft  
(0, 612)  
>>> first_page.mediaBox.upperRight  
(792, 612)
```

Each property returns a `tuple` containing the coordinates of the specified corner. You can access individual coordinates with square brackets, just like you would any other Python tuple:

```
>>> first_page.mediaBox.upperRight[0]  
792  
>>> first_page.mediaBox.upperRight[1]  
612
```

You can alter the coordinates of a `mediaBox` by assigning a new tuple to one of its properties:

```
>>> first_page.mediaBox.upperLeft = (0, 480)
>>> first_page.mediaBox.upperLeft
(0, 480)
```

When you change the `.upperLeft` coordinates, the `.upperRight` attribute adjusts automatically so that a rectangular shape is preserved:

```
>>> first_page.mediaBox.upperRight
(792, 480)
```

When you alter the coordinates of the `RectangleObject` returned by `.mediaBox`, you effectively crop the page. The `first_page` object now contains only the information present within the boundaries of the new `RectangleObject`.

Go ahead and write the cropped page to a new PDF file:

```
>>> pdf_writer = PdfFileWriter()
>>> pdf_writer.addPage(first_page)
>>> with Path("cropped_page.pdf").open(mode="wb") as output_file:
...     pdf_writer.write(output_file)
...
>>>
```

If you open the `cropped_page.pdf` file in your current working directory, you'll see that the top portion of the page has been removed.

How would you crop the page so that just the text on the left side of the page is visible? You need to cut the horizontal dimensions of the page in half. You can achieve this by altering the `.upperRight` coordinates of the `.mediaBox` object. Let's see how that works.

First you need to get new `PdfFileReader` and `PdfFileWriter` objects since we've just altered the first page in `pdf_reader` added it to '`pdf_writer`:

```
>>> pdf_reader = PdfFileReader(str(pdf_path))
>>> pdf_writer = PdfFileWriter()
```

Now get the fist page of the PDF:

```
>>> first_page = pdf_reader.getPage(0)
```

This time, let's work with a copy of the first page so that the page you just extracted stays intact. You can do that by importing the `copy` module from Python's standard library and using the `deepcopy()` function to make a copy of the page:

```
>>> import copy  
>>> left_side = copy.deepcopy(first_page)
```

Now we can alter `left_side` without changing the properties of `first_page`. That way we can use `first_page` later to extract the text on the right hand side of the page.

Now you need to do a little bit of math. We already worked out that we need to move the upper right hand corner of the `.mediaBox` to the top center of the page. To do that, you'll create a new tuple with the first component equal to half of the original value and assign it to the `.upperRight` property.

First, get the current coordinates of the upper right corner of the `.mediaBox`.

```
>>> current_coords = left_side.mediaBox.upperRight
```

Then create a new tuple whose first coordinate is half the value of the current coordinate, and second coordinate is the same as the original:

```
>>> new_coords = (current_coords[0] / 2, current_coords[1])
```

Finally, assign the new coordinates to the `.upperRight` property:

```
>>> left_side.mediaBox.upperRight = new_coords
```

You've now cropped the original page to contain only the text on the left side! Let's extract the right-hand side of the page next.

First get a new copy of `first_page`:

```
>>> right_side = copy.deepcopy(first_page)
```

To crop the page to just the right-hand side, move the `.upperLeft` corner instead of the `.upperRight` corner:

```
>>> right_side.mediaBox.upperLeft = new_coords
```

This sets the upper left corner to the same coordinates that you moved the upper right corner to when extracting the left-hand side of the page. So, `right_side.mediaBox` is now a rectangle whose upper left corner is at the top center of the page, and upper right corner is at the top right of the page.

Finally, add the `left_side` and `right_side` pages to `pdf_writer` and write them to a new PDF file:

```
>>> pdf_writer.addPage(left_side)
>>> pdf_writer.addPage(right_side)
>>> with Path("cropped_pages.pdf").open(mode="wb") as output_file:
...     pdf_writer.write(output_file)
...
>>>
```

Now open the `cropped_pages.pdf` file with a PDF reader. You should see a file with two pages, the first containing the text from the left-hand side of the original first page, and the second containing the text from the original right-hand side.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. In the Chapter 14 Practice Files folder there is a PDF called `split_and_rotate.pdf`. Create a new PDF called `rotated.pdf` in your home directory containing the pages of `split_and_rotate.pdf` rotated counterclockwise 90 degrees.
2. Using the `rotated.pdf` file you created in exercise 1, split each page

of the PDF vertically in the middle. Create a new PDF called `split.pdf` in your home directory containing all of the split pages.

`split.pdf` should have four pages with the numbers 1, 2, 3, and 4, in order.

[Leave feedback on this section »](#)

14.6 Encrypting and Decrypting PDFs

Sometimes PDF files are password protected. With the `PyPDF2` package, you can work with encrypted PDF files, as well as add password protection to existing PDFs.

PDF Encryption

The `.encrypt()` method of a `PdfFileWriter()` instance is used to add password protection to a PDF file. It has two main parameters:

1. `user_pwd` for setting the user password. This allows for opening and reading the PDF file.
2. `owner_pwd` for setting the owner password. This allows for opening the PDF without any restrictions, including editing.

Let's use `.encrypt()` to add a password to a PDF file. First, let's open the `newsletter.pdf` file in the Chapter 14 Practice Files directory:

```
>>> from pathlib import Path
>>> from PyPDF2 import PdfFileReader, PdfFileWriter
>>> pdf_path = (
...     Path.home() /
...     "python-basics-exercises" /
...     "ch13-interact-with-pdf-files" /
...     "practice_files" /
...     "newsletter.pdf"
... )
```

```
>>> pdf_reader = PdfFileReader(str(pdf_path))
```

Now create a new `PdfFileWriter` instance and add the pages from `pdf_reader` to it:

```
>>> pdf_writer = PdfFileWriter()
>>> pdf_writer.appendPagesFromReader(pdf_reader)
```

Now we'll add the password "SuperSecret" with the `pdf_writer.encrypt()` method:

```
>>> pdf_writer.encrypt(user_pwd="SuperSecret")
```

When you set only `user_pwd`, the `owner_pwd` argument defaults to the same string, so the above line of code sets both the user and owner passwords.

Finally write the encrypted PDF to an output file called `newsletter_protected.pdf` in your home directory:

```
>>> output_path = Path.home() / "newsletter_protected.pdf"
>>> with output_path.open(mode="wb") as output_file:
...     pdf_writer.write(output_file)
```

When you open the PDF with a PDF reader software you'll be prompted to enter a password. Enter "SuperSecret" to open the PDF.

If you need to set a separate owner password for the PDF, pass a second string to the `owner_pwd` parameter. For example, the following sets the user password to "SuperSecret" and the owner password to "ReallySuperSecret":

```
>>> pdf_writer.encrypt(user_pwd="SuperSecret", owner_pwd="ReallySuperSecret")
```

PDF Decryption

When you work with password-protected files programmatically, you need to decrypt them before you can access any of the contents.

To decrypt an encrypted PDF file, use the `.decrypt()` method of a `PdfFileReader` instance. The `.decrypt()` method has a single parameter called `password`. Let's open the encrypted `newsletter_protected.pdf` file you created in the previous section.

First, create a new `PdfFileReader` instance with the path to the protected PDF:

```
>>> from pathlib import Path  
>>> from PyPDF2 import PdfFileReader, PdfFileWriter  
>>> pdf_path = Path.home() / "newsletter_protected.pdf"  
>>> pdf_reader = PdfFileReader(str(pdf_path))
```

Before we decrypt the PDF, let's see what happens if we try to get the first page:

```
>>> pdf_reader.getPage(0)  
Traceback (most recent call last):  
  File "<stdin>", line 1, in <module>  
  File "/Users/damos/github/realpython/python-basics-exercises/venv/  
lib/python38-32/site-packages/PyPDF2/pdf.py", line 1176, in getPage  
    self._flatten()  
  File "/Users/damos/github/realpython/python-basics-exercises/venv/  
lib/python38-32/site-packages/PyPDF2/pdf.py", line 1505, in _flatten  
    catalog = self.trailer[ "/Root" ].getObject()  
  File "/Users/damos/github/realpython/python-basics-exercises/venv/  
lib/python38-32/site-packages/PyPDF2/generic.py",  
line 516, in __getitem__  
    return dict.__getitem__(self, key).getObject()  
  File "/Users/damos/github/realpython/python-basics-exercises/venv/  
lib/python38-32/site-packages/PyPDF2/generic.py",  
line 178, in getObject  
    return self.pdf.getObject(self).getObject()  
  File "/Users/damos/github/realpython/python-basics-exercises/venv/  
lib/python38-32/site-packages/PyPDF2/pdf.py", line 1617, in getObject  
    raise utils.PdfReadError("file has not been decrypted")  
PyPDF2.utils.PdfReadError: file has not been decrypted
```

A `PdfReadError` exception is raised that informs you that the PDF file has not been decrypted.

Go ahead and decrypt the file now:

```
>>> pdf_reader.decrypt(password="SuperSecret")  
1
```

`.decrypt()` returns an integer representing the success of the decryption:

1. 0 indicates the password is incorrect
2. 1 indicates the user password was matched
3. 2 indicated the owner password was matched

Once the file is decrypted you can access the contents of the PDF:

```
>>> pdf_reader.getPage(0)  
{'/Contents': IndirectObject(7, 0), '/CropBox': [0, 0, 612, 792],  
 '/MediaBox': [0, 0, 612, 792], '/Parent': IndirectObject(1, 0),  
 '/Resources': IndirectObject(8, 0), '/Rotate': 0, '/Type': '/Page'}
```

Now you can extract text, crop, or rotate pages to your heart's content!

Review Exercises

1. In the Chapter 14 Practice Files folder there is a PDF file called `top_secret.pdf`. Using `PdfFileWriter.encrypt()`, encrypt the file with the user password `Unguessable`.

Save the encrypted file as to your home directory with the filename `top_secret_encrypted.pdf`.

2. Open the `top_secret_encrypted.pdf` file you created in Exercise 1, decrypt it, and print the text contained on the first page of the PDF.

[Leave feedback on this section »](#)

14.7 Challenge: Unscramble A PDF

In the Chapter 14 Practice Files folder there is a PDF file called `scrambled.pdf` with seven pages. Each page contains a number 1 through 7, but they are out of order.

Additionally, some of the pages are rotated by one of 90, 180, or 270 degrees in either the clockwise or counterclockwise position.

Write a script that unscrambles the PDF by sorting the pages according to the number contained in the page text and rotating the page, if needed, so that it is upright.

Note

You can assume that every `PageObject` read from `scrambled.pdf` has a `"/Rotate"` key.

Save the unscrambled PDF to a file called `unscrambled.py` in your home directory.

You can find the solutions to this code challenge and many other bonus resources online at realpython.com/python-basics/resources.

[Leave feedback on this section »](#)

14.8 Create a PDF File From Scratch

The PyPDF2 package is great for reading and modifying existing PDF files, but it has a major limitation. You can't use it to create a new PDF file.

In this section, you'll use the `ReportLab` toolkit to generate PDF files from scratch. ReportLab is a full-featured PDF creation solution. There is a commercial version that costs money to use, but a limited-feature open source version is also available.

Note

This section is not meant to be an exhaustive introduction to ReportLab, but rather a sample of what is possible.

For more examples, checkout the ReportLab's [code snippet page](#).

Let's start by setting up your environment to work with ReportLab.

Install reportlab

To get started, you need to install ReportLab with pip:

```
$ python3 -m pip install reportlab
```

You can verify the installation with `pip show`:

```
$ python3 -m pip show reportlab
Name: reportlab
Version: 3.5.34
Summary: The Reportlab Toolkit
Home-page: http://www.reportlab.com/
Author: Andy Robinson, Robin Becker, the ReportLab team
        and the community
Author-email: reportlab-users@lists2.reportlab.com
License: BSD license (see license.txt for details),
        Copyright (c) 2000-2018, ReportLab Inc.
Location: c:/usersdaveavenv/lib/site-packages
Requires: pillow
Required-by:
```

Notice that the latest version of `reportlab` at the time of writing is 3.5.34. You'll need to restart IDLE if you have it open in order to use the `reportlab` package.

The Canvas Class

The main interface for creating PDFs with ReportLab is the `Canvas` class, which is located in the `reportlab.pdfgen.canvas` module.

Open a new IDLE interactive window and type the following to import the `Canvas` class:

```
>>> from reportlab.pdfgen.canvas import Canvas
```

When you make a new `Canvas` instance, you need to provide a string with the filename of the PDF you are creating. Let's create a new `Canvas` instance for the file `hello.pdf`:

```
>>> canvas = Canvas("hello.pdf")
```

You now have a `Canvas` instance associated with a file called `hello.pdf` in your current working directory assigned to the variable name `canvas`. The file `hello.pdf` does not exist yet, though.

Let's add some text to the PDF. To do that, you use the `.drawString()` method:

```
>>> canvas.drawString(72, 72, "Hello World")
```

The first two arguments passed to `.drawString()` determine the location on the canvas that the text is written. The first specifies **points** from the left edge of the canvas, and the second specifies points from the bottom edge.

A **point** is equal to 1/72 inches. So, 72 points is one inch, which means that the string "Hello World" is written one inch from the left and the bottom of the page.

To save the PDF to a file, use the `Canvas` object's `.save()` method:

```
>>> canvas.save()
```

You now have a PDF file called `hello.pdf` in your current working directory. You can open it with a PDF reader and see the text `Hello World`.

at the bottom of the page!

There are a few things to notice about the PDF you just created:

1. The default page size is A4, which is not the same as the American standard letter page size.
2. The font defaults to Helvetica with a default font size of 12 points.

You are not stuck with these settings.

Setting The Page Size

You can change the page size when you instantiate a `Canvas` object with the optional `pagesize` parameter. This parameter accepts a tuple of floating point values representing the width and height of the page in points.

For example, to set the page size to 8.5 inches width by 11 inches tall, you would create the following `Canvas`:

```
canvas = Canvas("hello.pdf", pagesize=(612.0, 729.0))
```

(612, 729) represents a letter sized paper because 8.5 times 72 is 612 and 11 times 72 is 729.

If doing the math to convert points to inches or centimeters isn't your cup of tea, you can use the `reportlab.lib.units` module to help you with the conversions. The `.units` module contains several helper objects, such as `inch` and `cm`, that simplify your conversions.

Go ahead and import the `inch` and `cm` objects from the `reportlab.lib.units` module:

```
>>> from reportlab.lib.units import inch, cm
```

Now let's inspect each object to see what they are:

```
>>> cm  
28.346456692913385  
>>> inch  
72.0
```

Both `cm` and `inch` are just floating point values. They represent the number of points contained in each unit. So `inch` is 72.0 points and `cm` is 28.346456692913385.

To use the units, multiply the unit name by the number of units you want to get the conversion to points. For example, here's how to use `inch` to set the page size to 8.5 inches wide by 11 inches tall:

```
>>> canvas = Canvas("hello.pdf", pagesize=(8.5 * inch, 11 * inch))
```

By passing a tuple to `pagesize`, you can create any size page that you want. However, the `reportlab` package has some standard page sizes built-in that are easier to work with.

The page sizes are located in the `reportlab.lib.pagesizes` module. For example, to set the page size to letter you can import the `LETTER` object from the `pagesizes` module and pass it to the `pagesize` parameter when instantiating your `Canvas`:

```
>>> from reportlab.lib.pagesizes import LETTER  
>>> canvas = Canvas("hello.pdf", pagesize=LETTER)
```

If you inspect the `LETTER` object, you'll see that it's a tuple of floats:

```
>>> LETTER  
(612.0, 792.0)
```

The `reportlab.lib.pagesize` module contains many standard page sizes. Here are some of them and their dimensions:

Page Size	Dimensions
A4	210 mm x 297 mm
LETTER	8.5 in x 11 in
LEGAL	8.5 in x 14 in

Page Size	Dimensions
TABLOID	11 in by 17 in

In addition to these, the module contains definitions for all of the [ISO 216](#) standard paper sizes.

Setting Font Properties

You can also change the font, font size, and font color when you write text to the canvas.

To change the font and font size, use the `Canvas.setFont()` method. First, create a new `Canvas` instance with file name `font-example.pdf` and a letter page size:

```
>>> canvas = Canvas("font-example.pdf", pagesize=LETTER)
```

Then set the font to Times New Roman with a size of 18 points:

```
>>> canvas.setFont("Times-Roman", 18)
```

Finally, write the string "Time New Roman (18 pt)" to the canvas and save it:

```
>>> canvas.drawString(1 * inch, 10 * inch, "Times New Roman (18 pt)")  
>>> canvas.save()
```

The text is written one inch from the left side of the page, and ten inches from the bottom. Open up the `font-example.pdf` file in your current working directory and check it out!

There are three fonts available by default:

1. "Courier"
2. "Helvetica"
3. "Times-Roman"

Each font has bold and italicized variants. Here's a list containing all of the font variations available in reportlab:

- "Courier"
- "Courier-Bold"
- "Courier-BoldOblique"
- "Courier-Oblique"
- "Helvetica"
- "Helvetica-Bold"
- "Helvetica-BoldOblique"
- "Helvetica-Oblique"
- "Times-Bold"
- "Times-BoldItalic"
- "Times-Italic"
- "Times-Roman"

You can also set the font color using the `Canvas.setFillColor()` method. In the following example, a PDF file named `font-colors.pdf` with blue text is created:

```
from reportlab.lib.colors import blue
from reportlab.lib.pagesizes import LETTER
from reportlab.lib.units import inch
from reportlab.pdfgen.canvas import Canvas

canvas = Canvas("font-colors.pdf", pagesize=LETTER)

# Set font to Times New Roman with 12 point size
canvas.setFont("Times-Roman", 12)

# Draw blue text one inch from the left and ten
# inches from the bottom
canvas.setFillColor("blue")
```

```
canvas.drawString(1*inch, 10*inch, "Black text")  
  
# Save the PDF file  
canvas.save()
```

Notice that the color `blue` is an object imported from the `reportlab.lib.colors` module. This module contains several common named colors. A full list of colors can be found in the [reportlab source code](#).

The examples in this section highlight the basics of working with ReportLab's `Canvas` object. But you've only scratched the surface. With ReportLab you can create tables, forms, and even high-quality graphics from scratch!

The [ReportLab User Guide](#) contains a plethora of examples of how to generate PDF documents from scratch. It's a great place to start if you're interested in learning more about creating PDFs with Python.

[Leave feedback on this section »](#)

14.9 Summary and Additional Resources

In this chapter, you learned how to create and modify PDF files with the `PyPDF2` and `reportlab` packages.

With `PyPDF2`, you learned how to:

- Read PDF files and extract text using the `PdfFileReader` class
- Write new PDF files using the `PdfFileWriter`
- Concatenate and merge PDF files using the `PdfFileMerger` class
- Rotate and Crop PDF pages
- Encrypt and decrypt PDF files with passwords

You also got an introduction to creating PDF files from scratch with the `reportlab` package. You learned about:

- The `Canvas` class
- Writing text to a `Canvas` with `.drawString()`
- Setting the font and font size with `.setFont()`
- Changing the font color with `.setFillColor()`

ReportLab is a powerful PDF creation tool, and you just scratched the surface of what's possible in this chapter.

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-13

Additional Resources

To learn more about working with PDF files in Python, check out the following resources:

- [How to Work with a PDF in Python](#)
- [ReportLab PDF Library User Guide](#)
- [Recommended resources on realpython.com](#)

[Leave feedback on this section »](#)

Chapter 15

Working With Databases

In Chapter 12 you learned how to store and retrieve data from files using Python. Another common way to store data is in a database.

A **database** is a structured system for storing data. It could be made up of several CSV files organized into directories, or something more elaborate.

Python comes with a light-weight SQL database called SQLite that is perfect for learning how to work with databases.

In this chapter, you will learn:

- How to create a SQLite database
- How to store and retrieve data from a SQLite database
- What packages are commonly used to work with other databases

Important

Some experience with SQL will be helpful when reading this chapter. If you want to learn more about SQL, check out the [resources on Real Python](#).

Let's dig in!

[Leave feedback on this section »](#)

15.1 An Introduction to SQLite

There are numerous SQL databases, and some are better suited to particular purposes than others. One of the simplest, most lightweight SQL databases is [SQLite](#), which runs directly on your machine and comes bundled with Python automatically.

In this section, you will learn how to use the `sqlite3` package to create a new database and store and retrieve data.

SQLite Basics

There are four basic steps to working with SQLite:

1. Import the `sqlite3` package
2. Connect to an existing database, or create a new database
3. Execute SQL statements on the database
4. Close the database connection

Let's get started by exploring these four steps in IDLEs interactive window. Open IDLE and type the following:

```
>>> import sqlite3  
>>> connection = sqlite3.connect("test_database.db")
```

The `sqlite3.connect()` function is used to connect to, or create, a database. When you execute `.connect("test_database.db")`, Python searches for an existing database called "test_database.db". If no database with that name is found, a new one is created in the current working directory. To create a database in a different directory, you must specify the full path in the argument to `.connect()`.

Note

If you want to create a one-time-use database while you're testing code or playing around with table structures, you can use the special name "`:memory:`" to create the database in temporary memory, like so:

```
connection = sqlite3.connect(":memory:")
```

The `.connect()` function returns a `sqlite3.Connection` object, which you can verify with the `type()` function:

```
>>> type(connection)
<class 'sqlite3.Connection'>
```

`Connection` objects represent the connection between your program and the database. They have several attributes and methods that can be used to interact with the database. To store and retrieve data, you need a `Cursor` object, which can be obtained with the `Connection.cursor()` function:

```
>>> cursor = connection.cursor()
>>> type(cursor)
<class 'sqlite3.Cursor'>
```

The `sqlite3.Cursor` object is your gateway to interacting with the database. Using a `Cursor`, you can create database tables, execute SQL statements, and fetch query results.

Note

The term **cursor** in database jargon usually refers to an object that is used to fetch results from a database query one row at a time. Although `sqlite3.Cursor` objects are used for this operation, they also do much more than is typically expected from a cursor. This is one important distinction to keep in mind when you use other databases besides SQLite.

Let's use the SQLite `datetime` function to get the current local time:

```
>>> query = "SELECT datetime('now', 'localtime');"  
>>> cursor.execute(query)  
<sqlite3.Cursor object at 0x000001A27EB85E30>
```

To get the current time, you first build a SQL statement with the correct syntax. In this case, `"SELECT datetime('now', 'localtime');"` is the statement we need, and it is assigned to the `query` variable. This returns the current time using the local time zone settings on your machine. Then the query is executed using the `cursor.execute()` method.

Note that `.execute()` returns a `Cursor` object, but we didn't assign this to a new variable. That's because `.execute()` alters the state of `cursor` and also returns the `cursor` object itself. This might look kind of strange, but it allows you to chain multiple `cursor` methods together on a single line.

You might be wondering where the time returned by the `datetime` function is. To get the query results, use the `cursor.fetchone()` method. `.fetchone()` returns a tuple containing the first row of results:

```
>>> cursor.fetchone()  
('2018-11-20 23:07:21',)
```

Since `.fetchone()` returns a tuple, you need to unpack the tuple elements to get the string containing the date and time information. Here's how you can do this by chaining the `.execute()` and `.fetchone()` methods:

```
>>> time = cursor.execute(query).fetchone()[0]  
>>> time  
'2018-11-20 23:09:45'
```

Finally, to close the database connection, use the `connection.close()` method:

```
>>> connection.close()
```

Using `with` to Manage Your Database Connection

Recall from Chapter 12 that you can use a `with` statement with the `open()` function to open the file and then automatically close the file once the `with` block has executed. The same pattern applies to SQLite database connections and is the recommended way to open a database connection.

Here's the `datetime` example from above using a `with` statement to manage the database connection:

```
>>> with sqlite3.connect("test_database.db") as connection:  
...     cursor = connection.cursor()  
...     query = "SELECT datetime('now', 'localtime');"  
...     time = cursor.execute(query).fetchone()[0]  
...  
>>> time  
'2018-11-20 23:14:37'
```

In this example, the `connection` variable is assigned to the `Connection` object returned by `sqlite3.connect()` in the `with` statement. The code in the `with` block gets a new `Cursor` object using `connection.cursor()`, and then gets the current time with the `Cursor` object's `.execute()` and `.fetchone()` methods.

Managing your database connections in a `with` statement has many advantages. The resulting code is often cleaner and shorter than code written without a `with` statement. Moreover, any changes made to the database are saved automatically, as you'll see in the next example.

Working With Database Tables

You don't usually want to create a whole database just to get the current time. Databases are used to store and retrieve information. To store data in a database, you need to create a table and write some values to it.

Let's create a table called `People` with three columns: `FirstName`,

`LastName`, and `Age`. The SQL query to create this table looks like this:

```
CREATE TABLE People(FirstName TEXT, LastName TEXT, Age INT);
```

Notice that `FirstName` and `LastName` have the word `TEXT` next to them, and `Age` is next to the word `INT`. This tells SQLite that values in the `FirstName` and `LastName` columns are text values, and values in the `Age` column are integers.

Once the table is created, you can populate it with some data using the `INSERT INTO` SQL command. The following query inserts the values `Ron`, `Obvious`, and `42` in the `FirstName`, `LastName`, and `Age` columns, respectively:

```
INSERT INTO People VALUES('Ron', 'Obvious', 42);
```

Note

Note that the string '`Ron`' and '`Obvious`' are delimited with single quotation marks. This makes them valid Python strings as well, but more importantly, only strings delimited with single quotes are valid SQLite strings.

When you write SQL queries as strings, you need to make sure that they are delimited with double quotation marks so that you can use single quotation marks inside of the Python strings to delimit SQLite strings.

SQLite is not the only SQL database that follows the single quote convention. Keep an eye out for this whenever you work with any SQL database.

Let's walk through how to execute these statements and save the changes to the database. First, we'll do it without using a `with` statement. Save and run the following script:

```
import sqlite3
```

```
connection = sqlite3.connect("test_database.db")
cursor = connection.cursor()
cursor.execute(
    """CREATE TABLE People(
        FirstName TEXT,
        LastName TEXT,
        Age INT
    );"""
)
cursor.execute(
    """INSERT INTO People VALUES(
        'Ron',
        'Obvious',
        42
    );"""
)
connection.commit()
connection.close()
```

First, you get a `Connection` object with `sqlite3.connect()` and assign it to the `connection` variable. A `Cursor` object is created with `connection.cursor()` and used to execute the two SQL statements for creating the `People` table and inserting some data.

The SQL statement in both `.execute()` methods have been written using triple quote strings so that we can format the SQL nicely. SQL ignores whitespace, so we can get away with this here and improve the readability of the Python code.

Finally, `connection.commit()` is used to save the data to the database. **Commit** is database jargon for saving data. If you do not run `connection.commit()`, no `People` table is created.

After the script runs, `test_database.db` has a `People` table with one row in it. You can verify this in the interactive window:

```
>>> connection = sqlite3.connect("test_database.db")
>>> cursor = connection.cursor()
>>> cursor.execute("SELECT * FROM People;")
<sqlite3.Cursor object at 0x000001F739DB6650>
>>> cursor.fetchone()
('Ron', 'Obvious', 42)
```

Next, let's look at the same script written using a `with` statement to manage the database connection. Before you can do anything, though, you need to delete the `People` table so that we can recreate it. Type the following into the interactive window to remove the `People` table from the database:

```
>>> cursor.execute("DROP TABLE People;")
<sqlite3.Cursor object at 0x000001F739DB6650>
>>> connection.commit()
>>> connection.close()
```

Now save and run the following script:

```
import sqlite3

with sqlite3.connect("test_database.db") as connection:
    cursor = connection.cursor()
    cursor.execute(
        """CREATE TABLE People(
            FirstName TEXT,
            LastName TEXT,
            Age INT
        );"""
    )
    cursor.execute(
        """INSERT INTO People VALUES(
            'Ron',
            'Obvious',
            42
        );"""
    )
```

```
)
```

Notice that not only is there no `connection.close()`, you also don't have to type `connection.commit()`. That's because any changes made to the database are automatically committed when the `with` block is done executing. This is another advantage to using a `with` statement to manage your database connection.

Executing Multiple SQL Statements

If you want to run more than one SQL statement at a time, you have a couple of options. One simple option is to use the `.executescript()` cursor method and give it a string that represents a full SQL script. Although semicolons separate lines of SQL code, it's common to pass a multiline string for readability. The following script does the same thing as the script you wrote at the beginning of this section:

```
import sqlite3

with sqlite3.connect("test_database.db") as connection:
    cursor = connection.cursor()
    cursor.executescript(
        """DROP TABLE IF EXISTS People;
        CREATE TABLE People(
            FirstName TEXT,
            LastName TEXT,
            Age INT
        );
        INSERT INTO People VALUES(
            'Ron',
            'Obvious',
            '42'
        );"""
    )
```

You can also execute many similar statements by using the `.executemany()` method and supplying a tuple of tuples, where

each inner tuple supplies the information for a single command. For instance, if you have a lot of people's information to insert into our People table, you can save this information in the following tuple of tuples:

```
people_values = (
    ("Ron", "Obvious", 42),
    ("Luigi", "Vercotti", 43),
    ("Arthur", "Belling", 28)
)
```

You can then insert all of these people at once in a single line of code:

```
cursor.executemany("INSERT INTO People VALUES(?, ?, ?)", people_values)
```

Here, the question marks act as place-holders for the tuples in `people_values`. This is called a **parameterized statement**. You may notice some similarity to this and formatting strings with the `.format()` string method you learned about in Chapter 4.

Avoid Security Issues With Parametrized Statements

For security reasons, especially when you need to interact with a SQL table based on the user input, you should *always* use parameterized SQL statements. This is because the user could potentially supply a value that looks like SQL code and causes your SQL statement to behave in unexpected ways. This is called a **SQL injection** attack and, even if you aren't dealing with a **malicious user**, it can happen entirely by accident.

For instance, suppose you want to insert a person into the `People` table based on user-supplied information. You might initially try something like the following:

```
import sqlite3
```

```
# Get person data from user
first_name = input("Enter your first name: ")
last_name = input("Enter your last name: ")
age = int(input("Enter your age: "))

# Execute insert statement for supplied person data
with sqlite3.connect("test_database.db") as connection:
    cursor = connection.cursor()
    cursor.execute(
        f"INSERT INTO People Values('{first_name}', '{last_name}', {age});"
    )
```

What if the user's name includes an apostrophe? Try adding Flannery O'Connor to the table, and you'll see that she breaks the code. This is because the apostrophe gets mixed up with the single quotes in the line, making it appear to the database that the SQL code ends earlier than expected.

In this case, the code only causes an error, which is bad enough. In some cases, though, bad input can corrupt an entire table. Many other hard-to-predict cases can break SQL tables, and even delete portions of your database. To avoid this, you should always use parameterized statements.

The following script does the same thing as the script above, but uses a parametrized statement to insert the user input into the database:

```
import sqlite3

first_name = input("Enter your first name: ")
last_name = input("Enter your last name: ")
age = int(input("Enter your age: "))
data = (first_name, last_name, age)

with sqlite3.connect("test_database.db") as connection:
    cursor = connection.cursor()
    cursor.execute("INSERT INTO People VALUES(?, ?, ?);", data)
```

You can update the content of a row by using a parametrized SQL UPDATE statement. For instance, if you want to change the Age associated with someone already in our `People` table, you could use the following:

```
cursor.execute(  
    "UPDATE People SET Age=? WHERE FirstName=? AND LastName=?;",  
    (45, 'Luigi', 'Vercotti')  
)
```

Retrieving Data

Of course, inserting and updating information in a database isn't all that helpful if you can't fetch that information from the database. To fetch data from a database, you can use the `.fetchone()` and `.fetchall()` cursor methods. These are similar to the `.readline()` and `.readlines()` methods for reading lines from a file. `.fetchone()` returns a single row from query results, while `.fetchall()` retrieves all of the results of a query at once.

The following script illustrates how to use `.fetchall()`:

```
import sqlite3  
  
values = (  
    ("Ron", "Obvious", 42),  
    ("Luigi", "Vercotti", 43),  
    ("Arthur", "Belling", 28),  
)  
  
with sqlite3.connect("test_database.db") as connection:  
    cursor = connection.cursor()  
    cursor.execute("DROP TABLE IF EXISTS People")  
    cursor.execute(  
        """CREATE TABLE People(  
            FirstName TEXT,  
            LastName TEXT,  
            Age INT
```

```

    );"""
)

cursor.executemany("INSERT INTO People VALUES(?, ?, ?);", values)

# Select all first and last names from people over age 30
cursor.execute(
    "SELECT FirstName, LastName FROM People WHERE Age > 30;"
)
for row in cursor.fetchall():
    print(row)

```

In the script above, you first drop the `People` table to destroy the changes made in the previous examples in this section. Then you create the `People` table and insert several values into it. Next, a `SELECT` statement is executed that returns the first and last names of all people over the age of 30.

Finally, `.fetchall()` returns the results of a query as a list of tuples, where each tuple contains the data from a single row in the query results. The output of the script looks like this:

```
('Ron', 'Obvious')
('Luigi', 'Vercotti')
```

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources. :

1. Create a new database with a table named `Roster` that has three fields: `Name`, `Species` and `IQ`. The `Name` and `Species` columns should be text fields, and the `IQ` column should be an integer field.
2. Populate your new table with the following values:

Name	Species	IQ
Jean-Baptiste Zorg	Human	122
Korben Dallas	Meat Popsicle	100

Name	Species	IQ
Ak'not	Mangalore	-5

3. Update the Species of Korben Dallas to be Human.
4. Display the names and IQs of everyone in the table classified as Human.

[Leave feedback on this section »](#)

15.2 Libraries for Working With Other SQL Databases

If you have a particular type of SQL database that you'd like to access through Python, most of the basic syntax is likely to be identical to what you just learned for SQLite. However, you'll need to install an additional package to interact with your database since SQLite is the only built-in option.

There are many SQL variants and corresponding Python packages available. A few of the most commonly used and reliable open-source alternatives to SQLite are:

- [pyodbc](#), which connects to ODBC (Open Database Connection) databases, such as Microsoft SQL Server
- [psycopg2](#), which connects to the PostgreSQL database
- [PyMySQL](#), which connects to MySQL databases

One difference between SQLite and other databases—besides the actual syntax of the SQL code, which changes slightly with most flavors of SQL—is that most databases require a username and password to connect. Check the documentation for the particular package you want to use to for the syntax for making a database connection.

The [SQLAlchemy](#) package is another popular option for working with databases. SQLAlchemy is an object-relational mapping, or ORM,

that uses an object-oriented paradigm to build database queries. It can be configured to connect to a variety of databases. The object-oriented approach allows you to make queries without writing and raw SQL statements.

[Leave feedback on this section »](#)

15.3 Summary and Additional Resources

In this chapter, you learned how to interact with the SQLite database that comes with Python. SQLite is a small and light SQL database that can be used to store and retrieve data in your Python programs. To interact with SQLite in Python, you must import the `sqlite3` module.

To work with an SQLite database, you first need to connect to existing database, or create a new database, with the `sqlite3.connect()` function, which returns a `Connection` object. Then you can use the `Connection.cursor()` method to get a new `Cursor` object.

`Cursor` objects are used to execute SQL statements and retrieve query results. For example, `Cursor.execute()` and `Cursor.executescript()` are used to execute SQL queries. You can retrieve query results using the `Cursor.fetchone()` and `Cursor.fetchall()` methods.

Finally, you learned about several third-party packages that you can use to connect to other SQL databases, including `psycopg2`, which is used to connect to PostgreSQL databases, and `pyodbc` for Microsoft SQL Server. You also learned about the `SQLAlchemy` library, which provides a standard interface for connecting to a variety of SQL databases.

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-14

Additional Resources

Here are some more resources on working with databases:

- [pyodbc Getting Started](#)
- [psycopg Documentation](#)
- [SQLAlchemy Tutorial](#)
- Recommended resources on realpython.com

[Leave feedback on this section »](#)

Chapter 16

Interacting With the Web

The Internet hosts perhaps the greatest source of information—and misinformation—on the planet.

Many disciplines, such as data science, business intelligence, and investigative reporting, can benefit enormously from collecting and analyzing data from websites.

Web scraping is the process of collecting and parsing raw data from the web, and the Python community has come up with some pretty powerful web scraping tools.

In this chapter, you will learn how to:

- Parse website data using string methods and regular expressions
- Parse website data using an HTML parser
- Interact with forms and other website components

Important

Some experience with **HTML**, short for **HyperText Markup Language**—will be helpful when reading this chapter. To learn more about HTML, check out the [resources on Real Python](#).

Let's go!

[Leave feedback on this section »](#)

16.1 Scrape and Parse Text From Websites

Collecting data from websites using an automated process is known as web scraping. Some websites explicitly forbid users from scraping their data with automated tools like the ones you will create in this chapter. Websites do this for either of two possible reasons:

1. The site has a good reason to protect its data. For instance, Google Maps doesn't let you to request too many results too quickly.
2. Making many repeated requests to a website's server may use up bandwidth, slowing down the website for other users and potentially overloading the server such that the website stops responding entirely.

Important

You should always check a website's acceptable use policy before scraping its data to see if accessing the website by using automated tools is a violation of its terms of use. Legally, web scraping against the wishes of a website is very much a [gray area](#).

Please be aware that the following techniques may be illegal when used on websites that prohibit web scraping.

Let's start by grabbing all of the HTML code from a single webpage. We'll take a [straightforward page](#) that's been set up just for practice:

```
from urllib.request import urlopen

url = "http://olympus.realpython.org/profiles/aphrodite"
html_page = urlopen(url)
```

```
html_text = html_page.read().decode("utf-8")  
  
print(html_text)
```

This displays the following result for us, which represents the full HTML of the page just as a web browser would see it:

```
<html>  
<head>  
<title>Profile: Aphrodite</title>  
</head>  
<body bgcolor="yellow">  
<center>  
<br><br>  
  
<h2>Name: Aphrodite</h2>  
<br><br>  
Favorite animal: Dove  
<br><br>  
Favorite color: Red  
<br><br>  
Hometown: Mount Olympus  
</center>  
</body>  
</html>
```

Calling `urlopen()` will cause the following error if Python cannot connect to the Internet:

```
URLError: <urlopen error [Errno 11001] getaddrinfo failed>
```

If you provide an invalid web address that can't be found, you will see the following error, which is equivalent to the "404" page that a browser would load:

```
HTTPError: HTTP Error 404: Not Found
```

Now we can scrape specific information from the webpage using `text`

parsing techniques. Text parsing involves looking through the full string of text and grabbing only the pieces that are relevant to us.

For instance, if we wanted to get the title of the webpage (in this case, “Profile: Aphrodite”), we could use the string `find()` method to search through the text of the HTML for the `<title>` tags and parse out the actual title using index numbers:

```
from urllib.request import urlopen

url = "http://olympus.realpython.org/profiles/aphrodite"
page = urlopen(url)
html = page.read().decode('utf-8')

start_tag = "<title>"
end_tag = "</title>"
start_index = html.find(start_tag) + len(start_tag)
end_index = html.find(end_tag)

print(html[start_index:end_index])
```

Running this script displays the HTML code limited to only the text in the title:

```
Profile: Aphrodite
```

Of course, this worked for a simple example, but HTML in the real world can be much more complicated and far less predictable. For a small taste of the “expectations versus reality” of text parsing, visit </profiles/poseidon> and view the HTML source code.

The HTML for the `/profiles/poseidon` page looks similar to the `/profiles/aphrodite` page, but there is a small difference. The opening `<title >` tag has an extra space in it before the closing `>` character. Re-run the script you used to parse the title from the `profiles/aphrodite` page, but this time set the `url` variable to `http://olympus.realpython.org/profiles/poseidon`.

Instead of just seeing the text `Profile: Poseidon`, you get the following:

```
<head>
<title >Profile: Poseidon
```

The modified script doesn't find the beginning of the `<title>` tag correctly because of that pesky space before the closing `>`. So, `html.find(end_tag)` returns `-1` because the exact string `<title>` wasn't found anywhere. When `-1` is added to `len(start_tag)`, which is `7`, the `start_index` variable gets assigned the value `6`.

The 6th character of the `html_text` string is the beginning `<` of the `<head>` tag. This means that `html[start_index:end_index]` returns all of the HTML starting with `<head>` and ending just before `</title>`.

These sorts of problems can occur in countless unpredictable ways. A more reliable alternative than using `find()` is to use **regular expressions**. **Regular expressions**—or “regex” for short—are strings that can be used to determine whether or not text matches a particular pattern.

Note

Regular expressions are not particular to Python. They are a general programming concept that can be used with a wide variety of programming languages. Regular expressions use a language all of their own that is notoriously difficult to learn but incredibly useful once mastered.

Python provides built-in support for regular expressions through the `re` module. Just as Python uses the backslash character as an “escape character” for representing special characters that can't simply be typed into strings, regular expressions use many different “special” characters (called **meta-characters**) that are interpreted as ways to signify different types of patterns.

For instance, the asterisk character, `*`, stands for “zero or more” of whatever came just before the asterisk. In the following example, the `re.findall()` function is used to find any text within a string

that matches a given regular expression. The first argument of `re.findall()` is the regular expression that you want to match, and the second argument is the string to test:

```
>>> import re

>>> re.findall("ab*c", "ac")
['ac']

>>> re.findall("ab*c", "abcd")
['abc']

>>> re.findall("ab*c", "acc")
['ac']

>>> re.findall("ab*c", "abcac")
['abc', 'ac']

>>> re.findall("ab*c", "abdc")
[]
```

Our regular expression, `ab*c`, matches any part of the string that begins with an “a,” ends with a “c,” and has zero or more of “b” in between the two. The `re.findall()` function returns a list of all matches. If no matches are found, an empty list is returned.

Note that the matching is case-sensitive. If you want to match this pattern regardless of upper-case or lower-case differences, you can pass a third argument with the value `re.IGNORECASE`, which is a specific variable stored in the `re` module:

```
>>> re.findall("ab*c", "ABC")
[]

>>> re.findall("ab*c", "ABC", re.IGNORECASE)
['ABC']
```

You can use a period `.` to stand for any single character in a regular

expression. For instance, we could find all the strings that contain the letters “a” and “c” separated by a single character as follows:

```
>>> re.findall("a.c", "abc")
['abc']

>>> re.findall("a.c", "abbc")
[]

>>> re.findall("a.c", "ac")
[]

>>> re.findall("a.c", "acc")
['acc']
```

Putting the term `.*` inside of a regular expression stands for any character repeated any number of times. For instance, `"a.*c"` can be used to find every substring that starts with “a” and ends with “c”, regardless of which letter—or letters—are in-between:

```
>>> re.findall("a.*c", "abc")
['abc']

>>> re.findall("a.*c", "abbc")
['abbc']

>>> re.findall("a.*c", "ac")
['ac']

>>> re.findall("a.*c", "acc")
['acc']
```

Often, you use the `re.search()` function to search for a particular pattern inside a string. This function is somewhat more complicated than `re.findall()` because it returns an object called a `MatchObject` that stores different “groups” of data. This is because there might be matches inside of other matches, and `re.search()` returns every possible result.

The details of the `MatchObject` object are irrelevant here. For now, just know that calling the `.group()` method on a `MatchObject` will return the first and most inclusive result, which in most instances is just what you want. For instance:

```
>>> match_results = re.search("ab*c", "ABC", re.IGNORECASE)
>>> match_results.group()
'ABC'
```

There is one more function in the `re` module that is useful for parsing out text. The `re.sub()` function, which is short for “substitute,” allows you to replace text in a string that matches a regular expression with new text (sort of like the `.replace()` method). The arguments passed to `re.sub()` are the regular expression, followed by the replacement text, followed by the string. For example:

```
>>> string = "Everything is <replaced> if it's in <tags>."
>>> string = re.sub("<.*>", "ELEPHANTS", string)
>>> string
'Everything is ELEPHANTS.'
```

Perhaps that wasn’t quite what you expected to happen.

The `re.sub()` function uses the regular expression "`<.*>`" to find and replace everything in between the first `<` and last `>`, which is most of the string. This is because Python’s regular expressions are **greedy**, meaning that they try to find the longest possible match when characters like `*` are used.

Alternatively, you can use the non-greedy matching pattern `*?`, which works the same way as `*` except that it matches the shortest possible string of text:

```
>>> string = "Everything is <replaced> if it's in <tags>."
>>> string = re.sub("<.*?>", "ELEPHANTS", string)
>>> string
"Everything is ELEPHANTS if it's in ELEPHANTS."
```

Armed with all this knowledge, let's now try to parse out the title from <http://olympus.realpython.org/profiles/dionysus>, which includes this rather carelessly written line of HTML:

```
<TITLE>Profile: Dionysus</title />
```

The `.find()` method would have a difficult time dealing with the inconsistencies here, but with the clever use of regular expressions, you can handle this code easily:

```
import re
from urllib.request import urlopen

url = "http://olympus.realpython.org/profiles/dionysus"
page = urlopen(url)
html = page.read().decode("utf-8")

pattern = "<title.*?>.*?</title.*?>"
match_results = re.search(pattern, html, re.IGNORECASE)
title = match_results.group()
title = re.sub("<.*?>", "", title) # Remove HTML tags

print(title)
```

Let's take a closer look at the first regular expression in the pattern string by breaking it down into three parts—`<title.*?>`, `.*?`, and `</title.*?>`.

1. `<title.*?>`—This pattern matches the opening `<TITLE>` tag in `html`. The `<title` part of the pattern matches with `<TITLE` because `re.search()` is called with `re.IGNORECASE`, and `.*?>` matches any text after `<TITLE` up to the first instance of `>`.
2. `.*?`—This pattern matches all text after the opening `<TITLE>` non-greedily, stopping at the first match for `</title.*?>`.
3. `</title.*?>`—The only difference between this pattern and the first one is the `/` character, so this matches the closing `</title />` tag in `html`.

The second regular expression, the string "<.*?>" also uses the non-greedy `.*?` to match all the HTML tags in the `title` string. By replacing any matches with "", the `re.sub()` function removes all of the tags returns only the text.

Regular expressions are a powerful tool when used correctly. This introduction barely scratches the surface. You can learn more about regular expressions and how to use them in the [Python Regular Expression HOWTO](#) section of the Python documentation.

Note

Web scraping can be tedious. No two websites are organized the same way, and HTML is often messy. Moreover, websites change over time. Web scrapers that work today are not guaranteed to work next year—or next week, for that matter!

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Write a script that grabs the full HTML from the page <http://olympus.realpython.org/profiles/dionysus>
2. Use the string `.find()` method to display the text following “Name:” and “Favorite Color:” (not including any leading spaces or trailing HTML tags that might appear on the same line).
3. Repeat the previous exercise using regular expressions. The end of each pattern should be a “<” (the start of an HTML tag) or a new-line character, and you should remove any extra spaces or newline characters from the resulting text using the string `.strip()` method.

[Leave feedback on this section »](#)

16.2 Use an HTML Parser to Scrape Websites

Although regular expressions are great for pattern matching in general, sometimes it's easier to use an HTML parser that is explicitly designed for parsing out HTML pages. There are many Python tools written for this purpose, but the [Beautiful Soup](#) library is a good one to start with.

To install Beautiful Soup, you can run the following in your terminal:

```
$ pip3 install beautifulsoup4
```

Run `pip show` to see the details of the package you just installed:

```
$ pip3 show beautifulsoup4
Name: beautifulsoup4
Version: 4.6.3
Summary: Screen-scraping library
Home-page: http://www.crummy.com/software/BeautifulSoup/bs4/
Author: Leonard Richardson
Author-email: leonardr@segfault.org
License: MIT
Location: c:\realpython\venv\lib\site-packages
Requires:
Required-by:
```

In particular, notice that the latest version at the time of writing is 4.6.3.

Once you have Beautiful Soup installed, you can now import the `bs4` module and pass a string of HTML to `BeautifulSoup` to begin parsing:

```
from bs4 import BeautifulSoup
from urllib.request import urlopen

url = "http://olympus.realpython.org/profiles/dionysus"
```

```
page = urlopen(url)
html = page.read().decode("utf-8")
soup = BeautifulSoup(html, "html.parser")
```

This script does three things:

1. The URL `http://olympus.realpython.org/profiles/dionysus` is opened using the `urlopen()` function from the `urllib.request` module.
2. The HTML from the page is read as a string and assigned to the `html` variable.
3. A `BeautifulSoup` object is created and assigned to the `soup` variable.

The `BeautifulSoup` object assigned to `soup` is created with two arguments. The first argument is the HTML to be parsed and the second argument, the string `"html.parser"`, tells the object which parser to use behind the scenes. `"html.parser"` represents Python's built-in HTML parser.

Save and run the above script in IDLE. When it is finished running, you can use the `soup` variable in the interactive window to parse the content of `html` in various ways.

For example, `BeautifulSoup` objects have a `.get_text()` method that can be used to extract all of the text from the document and remove any HTML tags automatically.

Type the following code into IDLE's interactive window:

```
>>> print(soup.get_text())
```

```
Profile: Dionysus
```

```
Name: Dionysus
```

```
Hometown: Mount Olympus
```

```
Favorite animal: Leopard
```

```
Favorite Color: Wine
```

There are a lot of blank lines in this output. These are the result of newline characters in the HTML document's text. You can remove these with the string `.replace()` method, if you need to.

Often, you only need to get specific text from an HTML document. Using Beautiful Soup to extract the text first and then using the `.find()` string method is *sometimes* easier than working with regular expressions.

However, sometimes the HTML tags themselves are the elements that point out the data you want to retrieve. For instance, perhaps you want to retrieve the URLs for all the images on the page. These links are contained in the `src` attribute of `` HTML tags. In this case, you can use the `find_all()` method to return a list of all instances of that particular tag:

```
>>> soup.find_all("img")
[, ]
```

This returns a list of all `` tags in the HTML document. The objects in the list look like they might be strings representing the tags, but they are actually instances of the `Tag` object provided by Beautiful Soup. `Tag` objects provide a simple interface for working with the information they contain.

Let's explore this a little by first unpacking the `Tag` objects from the list:

```
>>> image1, image2 = soup.find_all("img")
```

Each `Tag` object has a `.name` property that returns a string containing the HTML tag type:

```
>>> image1.name  
'img'
```

You can access the HTML attributes of the `Tag` object by putting their name in-between square brackets, just as if the attributes were keys in a dictionary.

For example, the `` tag has a single attribute `src` with the value `dionysus.jpg`. Likewise, and HTML tag such as the link `` has two attributes, `href` and `target`.

To get the source of the images in the Dionysus profile page, you access the `src` attribute using the dictionary notation mentioned above:

```
>>> image1["src"]  
'/static/dionysus.jpg'  
  
>>> image2["src"]  
'/static/grapes.png'
```

Certain tags in HTML documents can be accessed by properties of the `Tag` object. For example, to get the `<title>` tag in a document, you can use the `.title` property:

```
>>> soup.title  
<title>Profile: Dionysus</title>
```

If you look at the source of the Dionysus profile by navigating to the URL <http://olympus.realpython.org/profiles/dionysus>, right-clicking on the page, and selecting “View Page Source,” you will notice that the `<title>` tag as written in the document looks like this:

```
<title>Profile: Dionysus</title/>
```

Beautiful Soup automatically cleans up the tags for you by removing the extra space in the opening tag and the extraneous / in the closing tag.

You can also retrieve just the string in the title tag with the `.string` property of the `Tag` object:

```
>>> soup.title.string  
'Profile: Dionysus'
```

One of the more useful features of Beautiful Soup is the ability to search for specific kinds of tags whose attributes match certain values. For example, if we want to find all of the `` tags that have a `src` attribute equal to the value `/static/dionysus.jpg`, you can provide the following additional argument to the `.find_all()` method:

```
>>> soup.find_all("img", src="/static/dionysus.jpg")  
[]
```

This example is somewhat arbitrary, and the usefulness of this technique may not be apparent from the example. If you spend some time browsing various websites and viewing their page source, you'll notice that many websites have extremely complicated HTML structure.

When scraping data from websites, you are often interested in particular parts of the page. By spending some time looking through the HTML document, you can identify tags with unique attributes that can be used to extract the data you need.

Then, instead of relying on complicated regular expressions or using `.find()` to search through the document, you can directly access the particular tag you are interested in and extract the data you need.

In some cases, you may find that Beautiful Soup does not offer the functionality you need. The `lxml` library is somewhat trickier to get started with but offers far more flexibility than Beautiful Soup for pars-

ing HTML documents. You may want to check it out once you are comfortable with using Beautiful Soup.

Note

HTML parsers like Beautiful Soup can save you a lot of time and effort when it comes to locating specific data in webpages. However, sometimes HTML is so poorly written and disorganized that even a sophisticated parser like Beautiful Soup can't interpret the HTML tags properly.

In this case, you're often left to your own devices (namely, `.find()` and regex) to try to parse out the information you need.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Write a script that grabs the full HTML from the page <http://olympus.realpython.org/profiles>
2. Parse out a list of all the links on the page using Beautiful Soup by looking for HTML tags with the name `a` and retrieving the value taken on by the `href` attribute of each tag.
3. Get the HTML from each of the pages in the list by adding the full path to the file name, and display the text (without HTML tags) on each page using Beautiful Soup's `.get_text()` method.

[Leave feedback on this section »](#)

16.3 Interact With HTML Forms

The `urllib` module you have been working with so far this chapter is well suited for requesting the contents of a webpage. Sometimes, though, you need to interact with a webpage to obtain the content you

need. For example, you might need to submit a form or click on a button to display hidden content.

The Python standard library does not provide a built-in means for working with web pages interactively, but many third-party packages are available from PyPI. Among these, [MechanicalSoup](#) is a popular and relatively simple package to use.

In essence, Mechanical Soup installs what is known as a **headless browser**, which is a web browser with no graphical user interface. This browser is controlled programmatically via a Python script.

You can install Mechanical Soup with `pip3` in your terminal:

```
$ pip3 install MechanicalSoup
```

You can now view some details about the package with `pip3 show`:

```
$ pip3 show mechanicalsoup
Name: MechanicalSoup
Version: 0.10.0
Summary: A Python library for automating interaction with websites
Home-page: https://mechanicalsoup.readthedocs.io/
Author: UNKNOWN
Author-email: UNKNOWN
License: MIT
Location: c:\realpython\venv\lib\site-packages
Requires: requests, beautifulsoup4, six, lxml
Required-by:
```

In particular, notice that the latest version at the time of writing is 0.10.0.

Note

You may need to close and restart your IDLE session for MechanicalSoup to load and be recognized after it's been installed.

To get started, let's write a script that creates a new Browser instance

with Mechanical Soup and retrieves a webpage:

```
import mechanize

browser = mechanize.Browser()
url = "http://olympus.realpython.org/login"
page = browser.get(url)
```

If you save and run the above script, you can then access the `page` variable in IDLE's interactive window, which will be useful for following along with the rest of this section.

The `page` variable now stores various information returned by the web server. For example, you can access the HTML of the webpage through the `.soup` property:

```
>>> page.soup
```

This will print out the following HTML:

```
<html>
<head>
<title>Log In</title>
</head>
<body bgcolor="yellow">
<center>
<br/><br/>
<h2>Please log in to access Mount Olympus:</h2>
<br/><br/>
<form action="/login" method="post" name="login">
Username: <input name="user" type="text"/><br/>
Password: <input name="pwd" type="password"/><br/><br/>
<input type="submit" value="Submit"/>
</form>
</center>
</body>
</html>
```

The `/login` page accessed by the above script has a `<form>` on it with `<input>` elements for a username and password.

You should open this page in a browser and look at it yourself before moving on. Try typing in a random username and password combination. If you guessed incorrectly, the message “Wrong username or password!” is displayed at the bottom of the page.

However, if you provide the correct login credentials (username “zeus” and password “ThunderDude”), you are redirected to the `/profiles` page.

In the next example, you will see how to use Mechanical Soup to fill out and submit this form using Python!

The important section of HTML code is the login form—that is, everything inside the `<form>` tags. The `<form>` on this page has the `name` attribute set to `login`. This form contains two `<input>` elements, one named `user` and the other named `pwd`. The third `<input>` element is the “Submit” button.

Now that you know the underlying structure of the login form, as well as the credentials needed to log in, let’s take a look at a script that fills the form out and submits it:

```
import mechanize

# 1
browser = mechanize.Browser()
url = "http://olympus.realpython.org/login"
login_page = browser.get(url)
login_html = login_page.read()

# 2
form = login_html.select("form")[0]
form.select("input")[0]["value"] = "zeus"
form.select("input")[1]["value"] = "ThunderDude"
```

```
# 3  
profiles_page = browser.submit(form, login_page.url)
```

After saving and running the script, you can confirm that you successfully logged in by typing the following into the interactive window:

```
>>> profiles_page.url  
'http://olympus.realpython.org/profiles'
```

Let's break down the above example.

1. In the first part of the script, a `Browser` instance is created and used to request the `http://olympus.realpython.org/login` page. The HTML content of the page is assigned to the `login_html` variable using the `.soup` property.
2. The next section handles filling out the form. The first step is to retrieve the `<form>` element itself from the page's HTML. `login_html.select("form")` returns a list of all `<form>` elements on the page. Since the page has only one `<form>` element, you can access the form by retrieving the 0th element of the list. The next two lines select the username and password inputs and set their value to "zeus" and "ThunderDude", respectively.
3. Finally, the form is submitted with the `browser.submit()` method. Notice that two arguments are passed to this method, the `form` object and the URL of the `login_page`, which is accessed via `login_page.url`.

In the interactive window, you confirmed that the submission successfully redirected to the `/profiles` page. If something had gone wrong, the value of `profiles_page.url` would still be .

Note

We are always being encouraged to use long passwords with many different types of characters in them, and now you know the main reason: automated scripts like the one we just designed can be used by hackers to “brute force” logins by rapidly trying to log in with many different usernames and passwords until they find a working combination.

Besides this being highly illegal, almost all websites these days lock you out and report your IP address if they see you making too many failed requests, so don’t try it!

Now that we have the `profiles_page` variable set let’s see how to programmatically obtain the URL for each link on the `/profiles` page.

To do this, you use the `.select()` method again, this time passing the string “`a`” to select all of the `<a>` anchor elements on the page:

```
>>> links = profiles_page.soup.select("a")
```

Now you can iterate over each link and print the `href` attribute:

```
>>> for link in links:  
...     address = link["href"]  
...     text = link.text  
...     print(f"{text}: {address}")  
...  
Aphrodite: /profiles/aphrodite  
Poseidon: /profiles/poseidon  
Dionysus: /profiles/dionysus
```

The URLs contained in each `href` attribute are relative URLs, which aren’t very helpful if you want to navigate to them later using Mechanical Soup. If you happen to know the full URL, you can assign the portion needed to construct a full URL. In this case, the base URL is just <http://olympus.realpython.org>. Then you can concatenate the base URL with the relative URLs found in the `src` attribute:

```
>>> base_url = "http://olympus.realpython.org"
>>> for link in links:
...     address = base_url + link["href"]
...     text = link.text
...     print(f"{text}: {address}")
...
Aphrodite: http://olympus.realpython.org/profiles/aphrodite
Poseidon: http://olympus.realpython.org/profiles/poseidon
Dionysus: http://olympus.realpython.org/profiles/dionysus
```

You can do a lot with just the `.get()`, `.select()`, and `.submit()` methods. That said, Mechanical Soup's is capable of much more. To learn more about Mechanical Soup, check out the [official docs](#).

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Use Mechanical Soup to provide the correct username “zeus” and password “ThunderDude” to the login page submission form located at <http://olympus.realpython.org/login>.
2. Display the title of the current page to determine that you have been redirected to the `/profiles` page.
3. Use Mechanical Soup to return to [login page](#) by going “back” to the previous page.
4. Provide an incorrect username and password to the login form, then search the HTML of the returned webpage for the text “Wrong username or password!” to determine that the login process failed.

[Leave feedback on this section »](#)

16.4 Interact With Websites in Real-Time

Sometimes we want to be able to fetch real-time data from a website that offers continually updated information. In the dark days, before you learned Python programming, you would have been forced to sit in front of a browser, clicking the “Refresh” button to reload the page each time you want to check if updated content is available. Instead, you can easily automate this process using the `.get()` method of the Mechanical Soup Browser object.

Open up your browser of choice and navigate to <http://olympus.realpython.org/dice>. This page simulates a roll of a 6-sided die, updating the result each time you refresh the browser. As an example of working with real-time data, you will write a script that periodically scrapes this page for a new result. While this example is admittedly contrived, you will learn the basics of interacting with a website to retrieve periodically updated results.

The first thing you need to do is determine which element on the page contains the result of the die roll. Do this now by right-clicking anywhere on the page and clicking on “View page source.” A little more than halfway down the HTML code, there is an `<h2>` tag that looks like this:

```
<h2 id="result">4</h2>
```

The text of the `<h2>` tag might be different for you, but this is the page element you need to scrape the result.

Note

For this example, you can easily check that there is only one element on the page with `id="result"`. Although the `id` attribute is supposed to be unique, in practice you should always check that the element you are interested in is uniquely identified. If not, you need to be creative with how you select that element in your code.

Let's start by writing a simple script that opens the [/dice](#) page, scrapes the result, and prints it to the console:

```
import mechanicalsoup

browser = mechanicalsoup.Browser()
page = browser.get("http://olympus.realpython.org/dice")
tag = page.soup.select("#result")[0]
result = tag.text

print(f"The result of your dice roll is: {result}")
```

This example uses the BeautifulSoup `.select()` to find the element with `id=result`. The string "`#result`" passed to `.select()` uses the [CSS ID selector](#) `#` to indicate `result` is an `id` value.

To periodically get a new result, you'll need to create a loop that loads the page at each step of the loop. So everything below the line `browser = mechanicalsoup.Browser()` in the above script needs to go in the body of the loop.

For this example, let's get 4 rolls of the dice at 30-second intervals. To do that, the last line of your code needs to tell Python to pause running for 30 seconds. You can do this with the `sleep()` function from Python's `time` module. The `sleep()` function takes a single argument that represents the time to sleep in seconds. Here's a simple example to illustrate how the `sleep()` function works:

```
import time

print("I'm about to wait for five seconds...")
time.sleep(5)
print("Done waiting!")
```

If you run the above example, you see that the "Done waiting!" message isn't displayed until 5 seconds have passed since the first `print()` function is executed.

For the die roll example, you'll need to pass the number 30 to `sleep()`. Here's the updated script:

```
import time
import mechanicalsoup

browser = mechanicalsoup.Browser()

for i in range(4):
    page = browser.get("http://olympus.realpython.org/dice")
    tag = page.soup.select("#result")[0]
    result = tag.text
    print(f"The result of your dice roll is: {result}")
    time.sleep(30)
```

When you run the script, you will immediately see the first result printed to the console. After 30 seconds, the second result is displayed, then the third and finally the fourth. What happens after the fourth result is printed?

The script continues running for another 30 seconds before it finally stops!

Well, of *course* it does! That's what you told it to do! But it's kind of a waste of time. You can stop it from doing this by using an `if` statement to run the `time.sleep()` function only for the first three requests:

```
import time
import mechanicalsoup

browser = mechanicalsoup.Browser()

for i in range(4):
    page = browser.get("http://olympus.realpython.org/dice")
    tag = page.soup.select("#result")[0]
    result = tag.text
    print(f"The result of your dice roll is: {result}")

    # Wait 30 seconds if this isn't the last request
    if i < 3:
        time.sleep(30)
```

Note

With techniques like this, you can scrape data from websites that periodically update their data. However, you should be aware that requesting a page multiple times in rapid succession can be seen as suspicious, or even malicious, use of a website. It's possible to crash a server with an excessive volume of requests, so you can imagine that many websites are concerned about the volume of requests to their server!

Most websites publish a Terms of Use document. A link to this document can often be found in the website's footer. You should always read this document before attempting to scrape data from a website. If you can not find the Terms of Use, try to contact the website owner and ask them if they have any policies regarding request volume.

Failure to comply with the Terms of Use could result in your IP being blocked, so be careful and be respectful!

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Repeat the example in this section to scrape the die roll result, but additionally include the current time of the quote as obtained from the webpage. This time can be taken from part of a string inside a `<p>` tag that appears shortly after the result of the roll in the webpage's HTML.

[Leave feedback on this section »](#)

16.5 Summary and Additional Resources

Working with data from the Internet can be complicated. The structure of websites varies significantly from one site to the next, and even a single website can change often. Although it is possible to parse data from the web using tools in Python's standard library, there are many tools on PyPI that can help simplify the process.

In this chapter, you learned about Beautiful Soup and Mechanical Soup, two tools that help you write Python programs to automate website interactions. Beautiful Soup is used to parse HTML data collected from a website. Mechanical Soup is used to interact with website components, such as clicking on links and submitting forms. With tools like Beautiful Soup and Mechanical Soup, you can open up your programs to the world.

Web scraping techniques are used in many real-world disciplines. For example, investigative journalists rely on information collected from vast numbers of resources. Programmers have developed several tools for scraping, parsing, and processing data from websites to help journalists gather data and understand connections between people, places, and events.

Writing automated web scraping programs is fun. The Internet has no shortage of crazy content that can lead to all sorts of exciting projects. Just remember, not everyone wants you pulling data from their web servers. Always check a website's Terms of Use before you start scraping, and be respectful about how you time your web requests so that you don't flood a server with traffic.

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-15

Additional Resources

For more information on interacting with the web with Python, check out the following resources:

- [Practical Introduction to Web Scraping in Python](#)
- [API Integration in Python](#)
- Recommended resources on realpython.com

[Leave feedback on this section »](#)

Chapter 17

Scientific Computing and Graphing

Python is one of the leading programming languages in scientific computing and data science.

Python's popularity in this area is due, in part, to the wealth of third-party packages available on PyPI for manipulating and visualizing data.

From cleaning and manipulating large data sets, to visualizing data in plots and charts, Python's ecosystem has the tools you need to analyze and work with data.

In this chapter, you will learn how to:

- Work with arrays of data using `numpy`
- Create charts and plots with `matplotlib`

Let's dive in!

[Leave feedback on this section »](#)

17.1 Use NumPy for Matrix Manipulation

In this section, you will learn how to store and manipulate matrices of data using the [NumPy](#) package. Before getting to that, though, let's take a look at the problem NumPy solves.

If you have ever taken a course in linear algebra, you may recall that a matrix is a rectangular array of numbers. You can easily create a matrix in pure Python with a list of lists:

```
>>> matrix = [[1, 2, 3], [4, 5, 6], [7, 8, 9]]
```

This seemingly works well. You can access individual elements of the matrix using their indices. For example, to access the second element of the first row of the matrix, you would type:

```
>>> matrix[0][1]  
2
```

Now suppose you want to multiply every element of the matrix by 2. To do this, you need to write a nested `for` loop that loops over every element of each row of the matrix. You might use a nested `for` loop, like this:

```
>>> for row in matrix:  
...     for i in range(len(row)):  
...         row[i] = row[i] * 2  
...  
>>> matrix  
[[2, 4, 6], [8, 10, 12], [14, 16, 18]]
```

While this may not seem so hard, the point is that in pure Python, you need to do a lot of work from scratch to implement even simple linear algebra tasks. Think about what you need to do if you want to multiply two matrices together!

NumPy provides nearly all of the functionality you might ever need

out-of-the-box and is more efficient than pure Python. NumPy is written in the C language, and uses sophisticated algorithms for efficient computation, bringing you speed and flexibility.

Note

Even if you have no interest in using matrices for scientific computing, you still might find it helpful at some point to store data in a NumPy matrix because of the many useful methods and properties it provides.

For instance, perhaps you are designing a game and need an easy way to store, view and manipulate a grid of values with rows and columns. Rather than creating a list of lists or some other complicated structure, using a NumPy array is a simple way to store your two-dimensional data.

Install NumPy

Before you can work with NumPy, you'll need to install it using pip:

```
$ pip3 install numpy
```

Once NumPy has finished installing, you can see some details about the package by running `pip3 show`:

```
$ pip3 show numpy
Name: numpy
Version: 1.15.0
Summary: NumPy: array processing for numbers, strings,
         records, and objects.
Home-page: http://www.numpy.org
Author: Travis E. Oliphant et al.
Author-email: None
License: BSD
Location: c:\realpython\venv\lib\site-packages
Requires:
```

Required-by:

In particular, notice that the latest version at the time of writing is version 1.15.0.

Create a NumPy array

Now that you have NumPy installed let's create the same matrix from the first example in this section. Matrices in NumPy are instances of the `ndarray` object, which stands for “*n*-dimensional array.”

Note

An *n*-dimensional array is an array with *n* dimensions. For example, a 1-dimensional array is a list. A 2-dimensional array is a matrix. Arrays can also have 3, 4, or more dimensions.

In this section, we will focus on arrays with one or two dimensions.

To create an `ndarray` object, you can use the `array` alias. You initialize array objects with a list of lists, so to re-create the matrix from the first example as a NumPy array, you can do the following:

```
>>> import numpy as np  
>>> matrix = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])  
>>> matrix  
array([[1, 2, 3],  
       [4, 5, 6],  
       [7, 8, 9]])
```

Notice how NumPy displays the matrix in a conveniently readable format? This is even true when printing the matrix with the `print()` function:

```
>>> print(matrix)  
[[1 2 3]  
 [4 5 6]]
```

```
[7 8 9]]
```

Accessing individual elements of the array works just like accessing elements in a list of lists:

```
>>> matrix[0][1]  
2
```

You can optionally access elements with just a single set of square brackets by separating the indices with a comma:

```
>>> matrix[0, 1]  
2
```

At this point, you might be wondering what the major difference is between a NumPy array and a Python list. For starters, NumPy arrays can only hold objects of the same type (for instance, all numbers) whereas Python lists can hold mixed types of objects. Check out what happens if you try to create an array with mixed types:

```
>>> np.array([[1, 2, 3], ["a", "b", "c"]])  
array([[1, 2, 3],  
       ['a', 'b', 'c']], dtype='|<U11')
```

NumPy doesn't raise an error. Instead, the types are converted to match one another. In this case, NumPy converts every element to a string. The `dtype='|<U11'` that you see in the above output means that this array can only store Unicode strings whose length is at most 11 bytes.

On the one hand, the automatic conversions of data types can be helpful, but it can also be a potential source of frustration if the data types are not converted in the manner you expect. For this reason, it is generally a good idea to handle your type conversion *before* initializing an array object. That way you can be sure that the data type stored in your array matches your expectations.

Note

For more examples of how NumPy arrays differ from Python lists, checkout out this [FAQ answer](#).

In NumPy, each dimension in an array is called an `axis`. Both of the previous matrices you have seen have two axes. Arrays with two axes are also called **two-dimensional arrays**. Here is an example of a three-dimensional array:

```
>>> matrix = np.array([  
...     [[1, 2, 3], [4, 5, 6]],  
...     [[7, 8, 9], [10, 11, 12]],  
...     [[13, 14, 15], [16, 17, 18]]  
... ])
```

To access an element of the above array, you need to supply three indices:

```
>>> matrix[0][1][2]  
6  
  
>>> matrix[0, 1, 2]  
6
```

If you think creating the above three-dimensional array looks confusing, you'll see a better way to create higher dimensional arrays later in this section.

Array Operations

Once you have an `array` object created, you can start to unleash the power of NumPy and perform some operations.

Recall from the first example in this section how you had to write a nested `for` loop to multiply each element in a matrix by the number 2. In NumPy, this operation is as simple as multiplying your `array` object by 2:

```
>>> matrix = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
>>> 2 * matrix
array([[ 2,  4,  6],
       [ 8, 10, 12],
       [14, 16, 18]])
```

You can just as easily perform element-wise arithmetic on multi-dimensional arrays as well:

```
>>> second_matrix = np.array([[5, 4, 3], [7, 6, 5], [9, 8, 7]])
>>> second_matrix - matrix
array([[ 4,  2,  0],
       [ 3,  1, -1],
       [ 2,  0, -2]])
```

All of the basic arithmetic operators ($+$, $-$, $*$, $/$) operate on arrays element for element. For example, multiplying two arrays with the $*$ operator does *not* compute the product of two matrices. Consider the following example:

```
>>> matrix = np.array([[1, 1, 1], [1, 1, 1], [1, 1, 1]])
>>> matrix * matrix
array([[1, 1, 1],
       [1, 1, 1],
       [1, 1, 1]])
```

To calculate an actual [matrix product](#), you can use the `@` operator:

```
>>> matrix @ matrix
array([[3, 3, 3],
       [3, 3, 3],
       [3, 3, 3]])
```

Note

The @ operator was introduced in Python 3.5, so if you are using an older version of Python you must multiply matrices differently. NumPy provides a function called `matmul()` for multiplying two matrices:

```
>>> np.matmul(matrix, matrix)
array([[3, 3, 3],
       [3, 3, 3],
       [3, 3, 3]])
```

The @ operator actually relies on the `np.matmul()` function internally, so there is no real difference between the two methods.

Other common array operations are listed here:

```
>>> matrix = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])

>>> # Get a tuple of axis length
>>> matrix.shape
(3, 3)

>>> # Get an array of the diagonal entries
>>> matrix.diagonal()
array([1, 5, 9])

>>> # Get a 1-dimensional array of all entries
>>> matrix.flatten()
array([1, 2, 3, 4, 5, 6, 7, 8, 9])

>>> # Get the transpose of an array
>>> matrix.transpose()
array([[1, 4, 7],
       [2, 5, 8],
       [3, 6, 9]])
```

```
>>> # Calculate the minimum entry
>>> matrix.min()
1

>>> # Calculate the maximum entry
>>> matrix.max()
9

>>> # Calculate the average value of all entries
>>> matrix.mean()
5.0

>>> # Calculate the sum of all entries
>>> matrix.sum()
45
```

Stacking and Shaping Arrays

Two arrays can be stacked vertically using `np.vstack()` or horizontally using `np.hstack()` if their axis sizes match:

```
>>> A = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
>>> B = np.array([[10, 11, 12], [13, 14, 15], [16, 17, 18]])

>>> np.vstack([A, B])
array([[ 1,  2,  3],
       [ 4,  5,  6],
       [ 7,  8,  9],
       [10, 11, 12],
       [13, 14, 15],
       [16, 17, 18]])

>>> np.hstack([A, B])
array([[ 1,  2,  3, 10, 11, 12],
       [ 4,  5,  6, 13, 14, 15],
       [ 7,  8,  9, 16, 17, 18]])
```

You can also reshape arrays with the `np.reshape()` function:

```
>>> A.reshape(9, 1)
array([[1],
       [2],
       [3],
       [4],
       [5],
       [6],
       [7],
       [8],
       [9]])
```

Of course, the total size of the reshaped array must match the original array's size. For instance, you can't execute `matrix.reshape(2, 5)`:

```
>>> A.reshape(2, 5)
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
ValueError: cannot reshape array of size 9 into shape (2, 5)
```

In this case, you are trying to shape an array with 9 entries into an array with 2 columns and 5 rows. This requires a total of 10 entries.

The `np.reshape()` function can be particularly helpful in combination with `np.arange()`, which is NumPy's equivalent to Python's `range()` function. The main difference is that `np.arange()` returns an array object:

```
>>> matrix = np.arange(1, 10)
>>> matrix
array([1, 2, 3, 4, 5, 6, 7, 8, 9])
```

Just like with `range()`, `np.arange()` starts with the first argument and ends just before the second argument. So, `np.arange(1, 10)` returns an array containing the numbers 1 through 9.

Together, `np.arange()` and `np.reshape()` provide a useful way to create

a matrix:

```
>>> matrix = matrix.reshape(3, 3)
>>> matrix
array([[1, 2, 3],
       [4, 5, 6],
       [7, 8, 9]])
```

You can even do this in a single line by chaining the calls to `np.arange()` and `np.reshape()` together:

```
>>> np.arange(1, 10).reshape(3, 3)
array([[1, 2, 3],
       [4, 5, 6],
       [7, 8, 9]])
```

This technique for creating matrices is particularly useful for creating higher-dimensional arrays. Here's how to create a three-dimensional array using `np.array()` and `np.reshape()`:

```
>>> np.arange(1, 13).reshape(3, 2, 2)
array([[[ 1,  2],
        [ 3,  4]],
       [[ 5,  6],
        [ 7,  8]],
       [[ 9, 10],
        [11, 12]]])
```

Of course, not every multi-dimensional array can be built from a sequential list of numbers. In that case, it is often easier to create and flat, one-dimensional list of entries and then `np.reshape()` the array into the desired shape:

```
>>> arr = np.array([1, 3, 5, 7, 9, 11, 13, 15, 17, 19, 21, 23])
>>> arr.reshape(3, 2, 2)
```

```
array([[[ 1,  3],  
       [ 5,  7]],  
  
      [[ 9, 11],  
       [13, 15]],  
  
      [[17, 19],  
       [21, 23]]])
```

In the list passed to `np.array()` in the above example, the difference between any pair of consecutive numbers is 2. You can simplify the creation of these kinds of arrays by passing an optional third argument to the `np.arange()` called the **stride**:

```
>>> np.arange(1, 24, 2)  
array([ 1,  3,  5,  7,  9, 11, 13, 15, 17, 19, 21, 23])
```

With that in mind, you can re-write the previous example even more simply:

```
>>> np.arange(1, 24, 2).reshape(3, 2, 2)  
array([[[ 1,  3],  
       [ 5,  7]],  
  
      [[ 9, 11],  
       [13, 15]],  
  
      [[17, 19],  
       [21, 23]]])
```

Sometimes you need to work with matrices of random data. With NumPy, creating random matrices is easy. The following creates a random 3×3 matrix:

```
>>> np.random.random([3, 3])  
array([[0.27721176, 0.66206403, 0.20722988],  
      [0.15722803, 0.06286636, 0.47220672],
```

```
[0.55657541, 0.27040345, 0.24558674]])
```

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Create a 3×3 NumPy array named `first_matrix` that includes the numbers 3 through 11 by using `np.arange()` and `np.reshape()`.
2. Display the minimum, maximum and mean of all entries in `first_matrix`.
3. Square every entry in `first_matrix` using the `**` operator, and save the results in an array named `second_matrix`.
4. Use `np.vstack()` to stack `first_matrix` on top of `second_matrix` and save the results in an array named `third_matrix`.
5. Use the `@` operator to calculate the matrix product of `third_matrix` by `first_matrix`.
6. Reshape `third_matrix` into an array of dimensions $3 \times 3 \times 2$.

[Leave feedback on this section »](#)

17.2 Use `matplotlib` for Plotting Graphs

In the previous section, you learned how to work with arrays of data using the `NumPy` package. While `NumPy` makes working with and manipulating data simple, it does not provide a means for human consumption of data. For that, you need to visualize your data.

Data visualization is a broad topic, complete with its own theory and a host of tools for displaying and interacting with visualizations. In this section, you will get an introduction to the `matplotlib` package, which is one of the more popular packages for quickly creating two-dimensional figures. Initially released in 2003, `matplotlib` is one of the oldest Python plotting libraries available. It remains popular and is still being actively developed to this day.

If you have ever created graphs in MATLAB, you will find that `matplotlib` in many ways directly emulates this experience. The similarities between MATLAB and `matplotlib` are intentional. The MATLAB plotting interface was a direct inspiration for `matplotlib`. Even if you haven't used MATLAB, you will likely find creating plots with `matplotlib` to be simple and straightforward.

Let's dive in!

Install `matplotlib`

You can install `matplotlib` from your terminal with `pip3`:

```
pip3 install matplotlib
```

You can then view some details about the package with `pip3 show`:

```
$ pip3 show matplotlib
Name: matplotlib
Version: 2.2.3
Summary: Python plotting package
Home-page: http://matplotlib.org
Author: John D. Hunter, Michael Droettboom
Author-email: matplotlib-users@python.org
License: BSD
Location: c:\realpython\venv\lib\site-packages
Requires: python-dateutil, pytz, kiwisolver, numpy,
          cycler, six, pyparsing
Required-by:
```

In particular, note that the latest version at the time of writing is version 2.2.3.

Basic Plotting With `pyplot`

The `matplotlib` package provides two distinct means of creating plots. The first, and simplest, method is through the `pyplot` interface. This is the interface that MATLAB users will find the most familiar.

The second method for plotting in `matplotlib` is through what is known as the [object oriented API](#). The object-oriented approach offers more control over your plots than is available through the `pyplot` interface. However, the concepts are generally more abstract.

In this section, you'll learn how to get up and running with the `pyplot` interface. You'll be pumping out some great looking plots in no time!

Note

The developers of `matplotlib` suggest you try to use the object-oriented API instead of the `pyplot` interface. In practice, if the `pyplot` interface offers you everything you need, then don't be ashamed to stick with it!

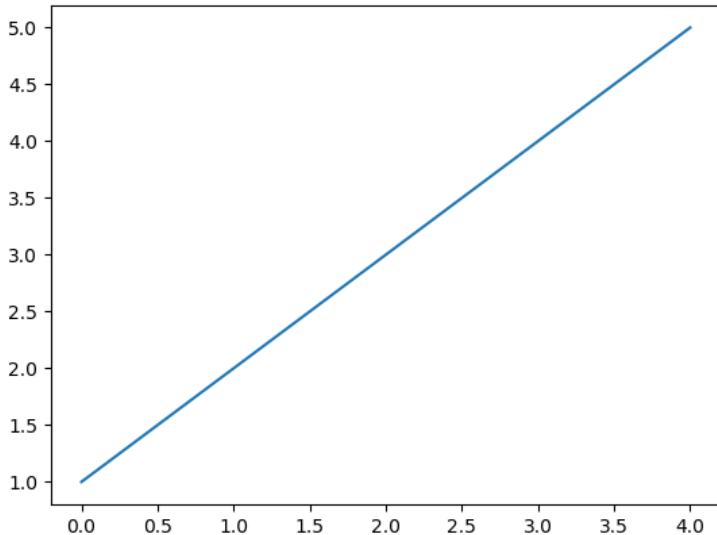
That said, if you are interested in learning more about the object-oriented approach, check out Real Python's [Python Plotting With Matplotlib \(Guide\)](#).

Let's start by creating a simple plot. Open IDLE and run the following script:

```
from matplotlib import pyplot as plt

plt.plot([1, 2, 3, 4, 5])
plt.show()
```

A new window appears displaying the following plot:



In this simple script, you created a plot with just a single line of code. The line `plt.plot([1, 2, 3, 4, 5])` creates a plot with a line through the points (0, 1), (1, 2), (2, 3), (3, 4), and (4, 5). The list [1, 2, 3, 4, 5] that you passed to the `plt.plot()` function represents the y-values of the points in the plot. Since you didn't specify any x-values, `matplotlib` automatically uses the indices of the list elements which, since Python starts counting at 0, are 0, 1, 2, 3 and 4.

The `plt.plot()` function creates a plot, but it does not display anything. The `plot.show()` function must be called to display the plot.

Note

If you are working in Windows, you should have no problem recreating the above plot from IDLE's interactive window. However, some operating systems have trouble displaying plots with `plot.show()` when called from the interactive window. We recommend working through each example in a new script.

If `plt.show()` works from the interactive window on your machine and you decide to follow along that way, be aware that once the figure is displayed in the new window, control isn't returned to the interactive window until you close the figure's window. That is, you won't see a new `>>>` prompt until the figure's window has been closed.

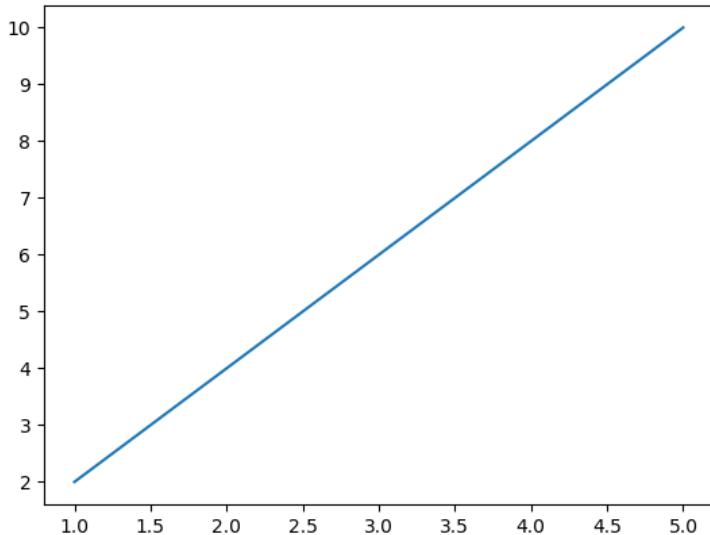
You can specify the x-values for the points in your plot by passing two lists to the `plt.plot()` function. When two arguments are provided to `plt.plot()`, the first list specifies the x-values and the second list specifies the y-values:

```
from matplotlib import pyplot as plt

xs = [1, 2, 3, 4, 5]
ys = [2, 4, 6, 8, 10]

plt.plot(xs, ys)
plt.show()
```

Running the above script produces the following plot:



At first glance, this figure may look exactly like the first. However, the labels on the axes now reflect the new x- and y-coordinates of the points.

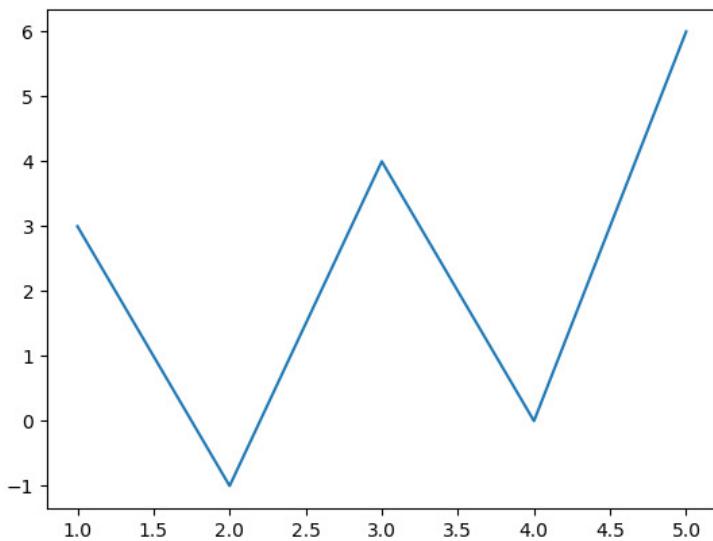
You can use `plot()` to plot more than lines. In the graphs above, the points being plotted just happen to all fall on the same line. By default, when plotting points with `.plot()`, each pair of consecutive points being plotted is connected with a line segment.

The following plot displays some data that doesn't fall on a line:

```
from matplotlib import pyplot as plt

xs = [1, 2, 3, 4, 5]
ys = [3, -1, 4, 0, 6]

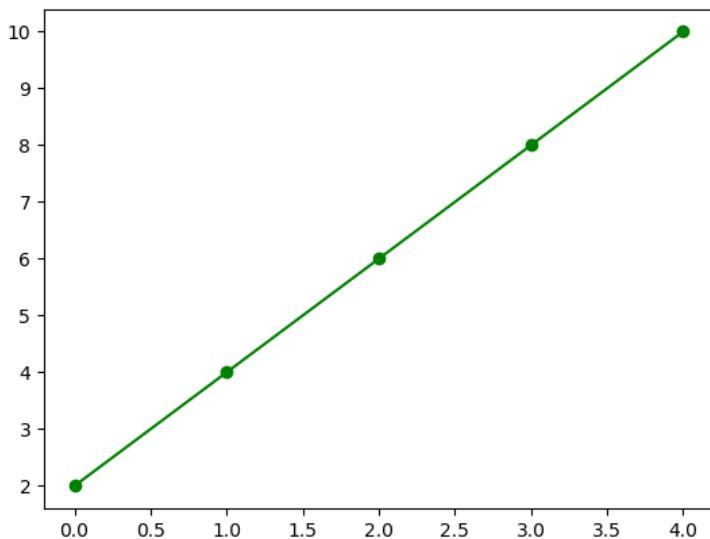
plt.plot(xs, ys)
plt.show()
```



There is an optional “formatting” argument that can be inserted into `plot()` after specifying the points to be plotted. This argument specifies the color and style of lines or points to draw.

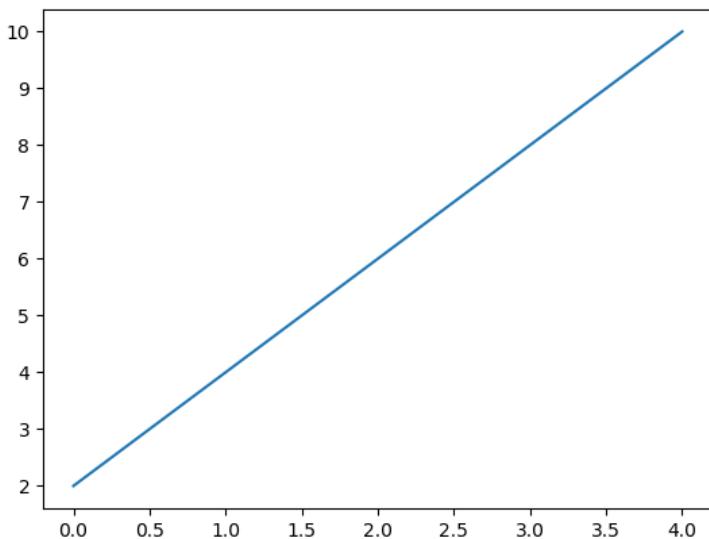
Unfortunately, the standard is borrowed from MATLAB and (compared to most Python) the formatting is not very intuitive to read or remember. The default value is “solid blue line,” which would be represented by the format string `b-`. If we wanted to plot green circular dots connected by solid lines instead, we would use the format string `g-o` like so:

```
from matplotlib import pyplot as plt  
  
plt.plot([2, 4, 6, 8, 10], "g-o")  
plt.show()
```



Note

For reference, the full list of possible formatting combinations can be found [here](#).



Plot Multiple Graphs in the Same Window

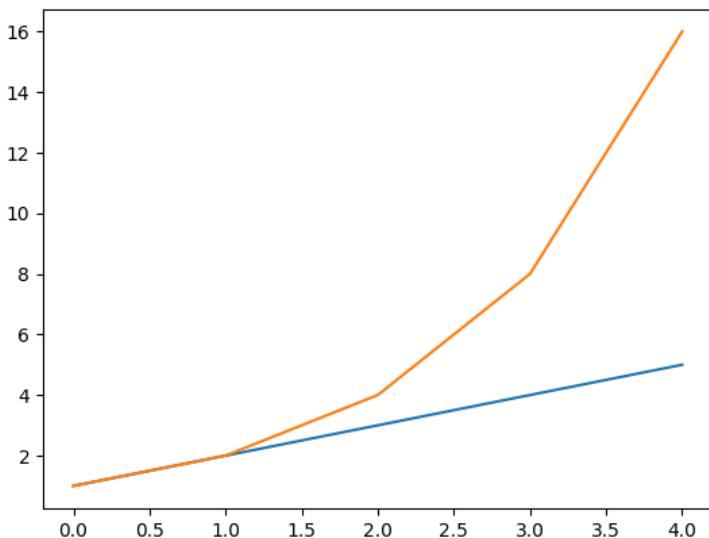
If you need to plot multiple graphs in the same window, you can do so a few different ways.

You can pass multiple pairs of x- and y-value lists:

```
from matplotlib import pyplot as plt

xs = [0, 1, 2, 3, 4]
y1 = [1, 2, 3, 4, 5]
y2 = [1, 2, 4, 8, 16]

plt.plot(xs, y1, xs, y2)
plt.show()
```



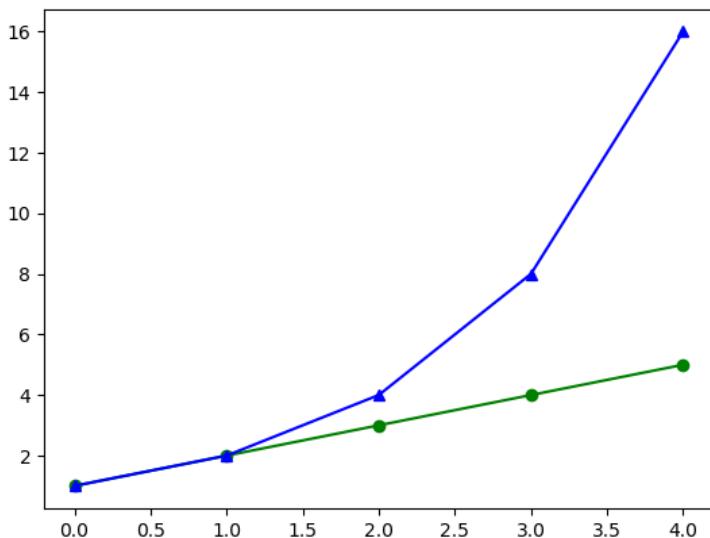
Notice that each graph is displayed in a different color. This built-in functionality of the `plot()` function is convenient for making easy-to-read plots very quickly.

If you want to control the style of each graph, you can pass the formatting strings to the `plot()` in addition to the x- and y-values:

```
from matplotlib import pyplot as plt

xs = [0, 1, 2, 3, 4]
y1 = [1, 2, 3, 4, 5]
y2 = [1, 2, 4, 8, 16]

plt.plot(xs, y1, "g-o", xs, y2, "b-^")
plt.show()
```



Passing multiple sets of points to `plot()` may work well when you only have a couple of graphs to display, but if you need to show many, it might make more sense to display each one with its own `plot()` function.

For example, the following script displays the same plot as the previous example:

```
from matplotlib import pyplot as plt

plt.plot([1, 2, 3, 4, 5], "g-o")
plt.plot([1, 2, 4, 8, 16], "b-^")
plt.show()
```

Plot Data From NumPy Arrays

Up to this point, you have been storing your data points in pure Python lists. In the real world, you will most likely be using some-

thing like a NumPy array to store your data. Fortunately, `matplotlib` plays nicely with array objects.

Note

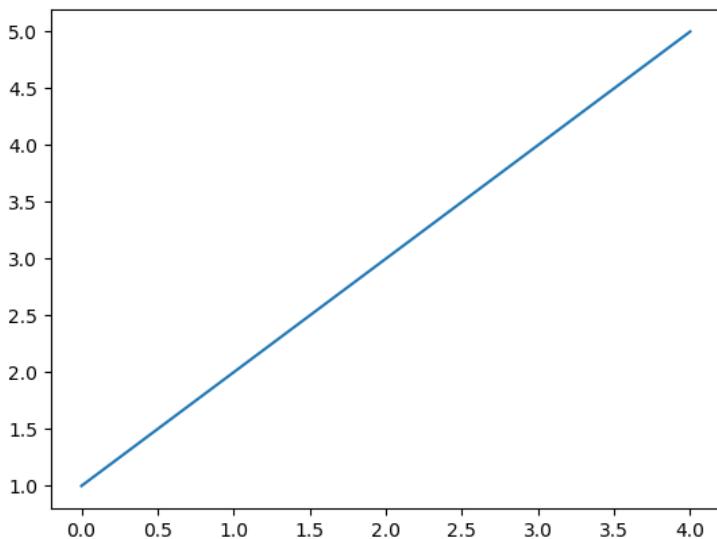
If you do not currently have NumPy installed, you need to install it with `pip`. For more information, please refer to the previous section in this chapter.

For example, instead of a `list`, you can use NumPy's `arange()` function to define your data points and then pass the resulting `array` object to the `plot()` function:

```
from matplotlib import pyplot as plt
import numpy as np

array = np.arange(1, 6)

plt.plot(array)
plt.show()
```



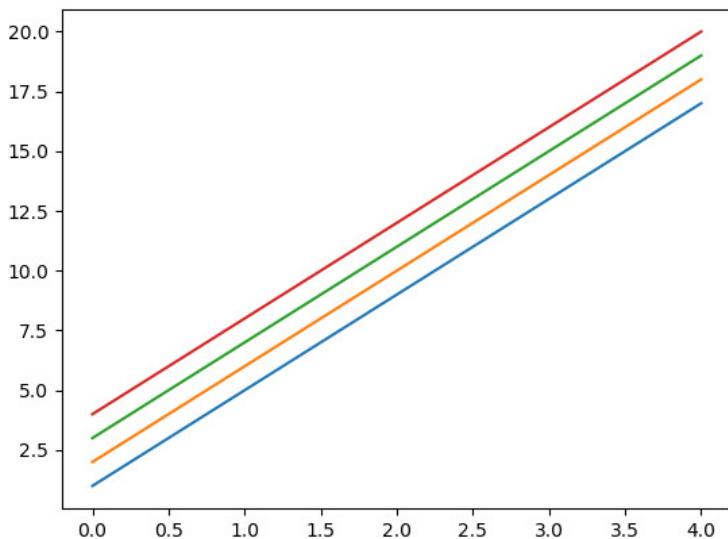
Passing a two-dimensional array plots each *column* of the array as the y-values for a graph. For example, the following script plots four lines:

```
from matplotlib import pyplot as plt
import numpy as np

data = np.arange(1, 21).reshape(5, 4)

# data now contains the following array:
# array([[ 1,  2,  3,  4],
#        [ 5,  6,  7,  8],
#        [ 9, 10, 11, 12],
#        [13, 14, 15, 16],
#        [17, 18, 19, 20]])

plt.plot(data)
plt.show()
```

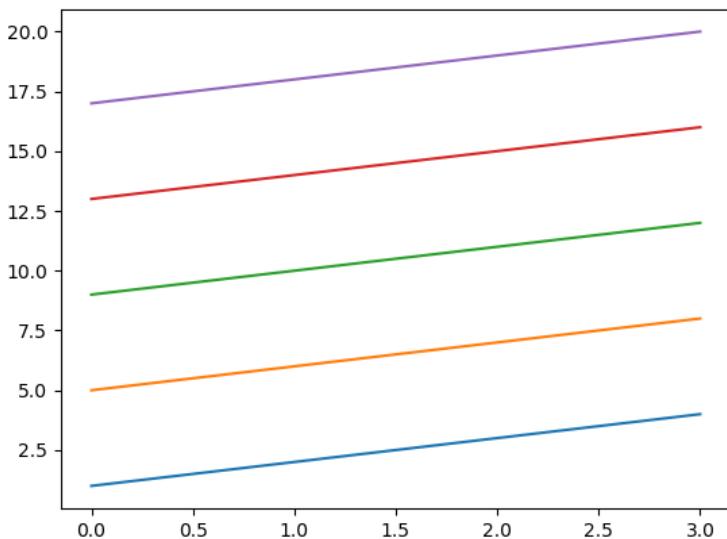


If instead you want to plot the rows of the matrix, you need to plot the transpose of the array. The following script plots the five rows of the same array from the previous example:

```
from matplotlib import pyplot as plt
import numpy as np

data = np.arange(1, 21).reshape(5, 4)

plt.plot(data.transpose())
plt.show()
```



Format Your Plots to Perfection

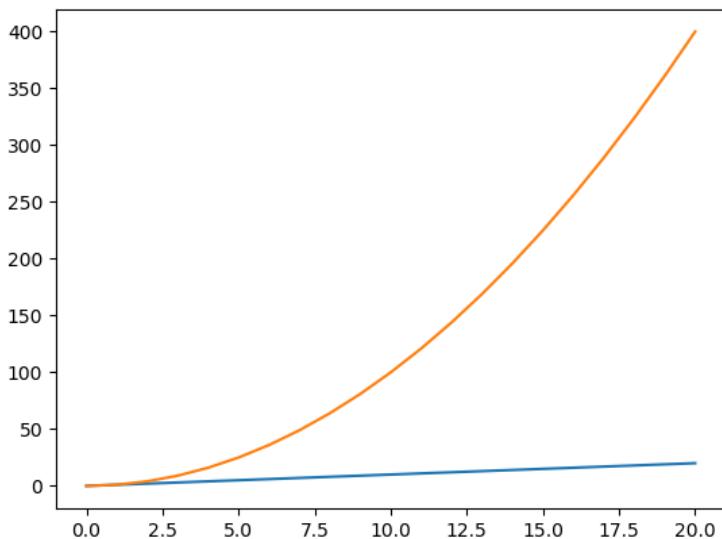
So far, the plots you have seen don't provide any information about what the plot represents. In this section, you will learn how to change the format and layout of your plots to make them easier to understand.

Let's start by plotting the amount of Python learned in the first 20 days of reading Real Python versus another website:

```
from matplotlib import pyplot as plt
import numpy as np

days = np.arange(0, 21)
other_site = np.arange(0, 21)
real_python = other_site ** 2

plt.plot(days, other_site)
plt.plot(days, real_python)
plt.show()
```



As you can see, the gains from reading Real Python are exponential! However, if you showed this graph to someone else, they may not understand what's going on.

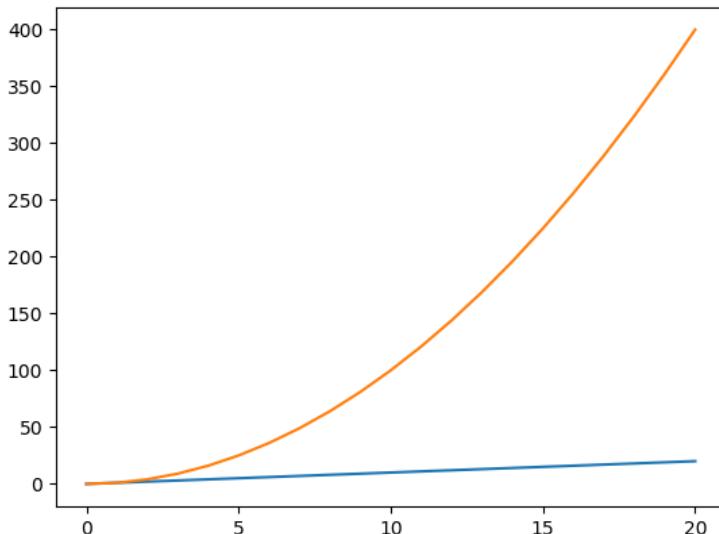
First of all, the x-axis is a little weird. It is supposed to represent days but is displaying half days instead. It would also be helpful to know what each line and axis represents. A title describing the plot wouldn't hurt, either.

Let's start with adjusting the x-axis. You can use the `plt.xticks()` function to specify where the ticks should be located by passing a list of locations. If we pass the list [0, 5, 10, 15, 20], the ticks should mark days 0, 5, 10, 15 and 20:

```
from matplotlib import pyplot as plt
import numpy as np
```

```
days = np.arange(0, 21)
other_site = np.arange(0, 21)
real_python = other_site ** 2

plt.plot(days, other_site)
plt.plot(days, real_python)
plt.xticks([0, 5, 10, 15, 20])
plt.show()
```



Nice! That's a little easier to read, but it still isn't clear what each axis represents.

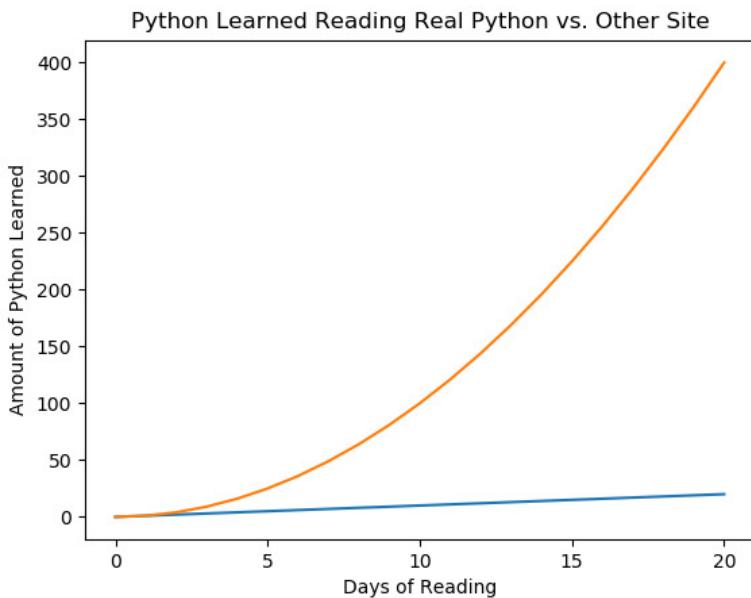
You can use the `plt.xlabel()` and `plt.ylabel()` to label the x- and y-axes, respectively. Just provide a string as an argument, and `matplotlib` displays the label on the corresponding axis.

While we're labeling things, let's go ahead and give the plot a title with the `plt.title()` function:

```
from matplotlib import pyplot as plt
import numpy as np

days = np.arange(0, 21)
other_site = np.arange(0, 21)
real_python = other_site ** 2

plt.plot(days, other_site)
plt.plot(days, real_python)
plt.xticks([0, 5, 10, 15, 20])
plt.xlabel("Days of Reading")
plt.ylabel("Amount of Python Learned")
plt.title("Python Learned Reading Real Python vs Other Site")
plt.show()
```



Now we're starting to get somewhere!

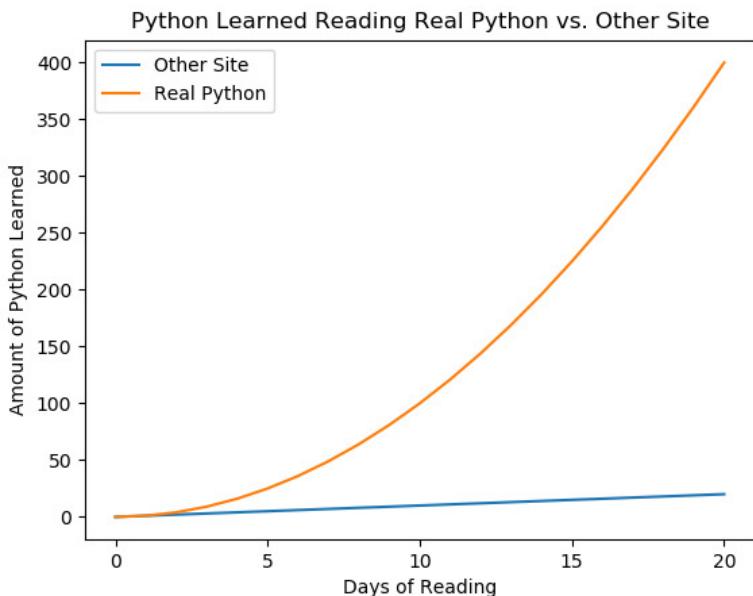
There's only one problem. It's not clear which graph represents Real Python and which one represents the other website.

To clarify which graph is which, you can add a legend with the `plt.legend()` function. The primary argument of the `legend()` function is a list of strings identifying each graph in the plot. These strings must be ordered in the same order the graphs were added to the plot:

```
from matplotlib import pyplot as plt
import numpy as np

days = np.arange(0, 21)
other_site = np.arange(0, 21)
real_python = other_site ** 2

plt.plot(days, other_site)
plt.plot(days, real_python)
plt.xticks([0, 5, 10, 15, 20])
plt.xlabel("Days of Reading")
plt.ylabel("Amount of Python Learned")
plt.title("Python Learned Reading Real Python vs Other Site")
plt.legend(["Other Site", "Real Python"])
plt.show()
```

**Note**

There are many ways to customize legends. For more information, check out the [Legend Guide](#) in the `matplotlib` documentation.

Other Types of Plots

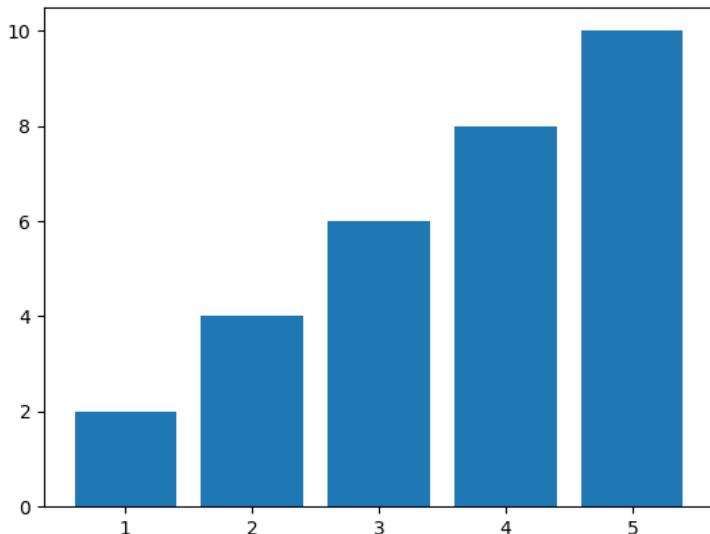
Aside from line charts, which up until now you have seen exclusively, `matplotlib` provides simple methods for creating other kinds of charts.

One frequently used type of plot in basic data visualization is the bar chart. You can easily create bar charts using the `plt.bar()` function. You must provide at least two arguments to `bar()`. The first is a list of x-values for the center point for each bar, and the second is the value for the top of each bar:

```
from matplotlib import pyplot as plt

xs = [1, 2, 3, 4, 5]
tops = [2, 4, 6, 8, 10]

plt.bar(xs, tops)
plt.show()
```



Just like the `plot()` function, you can use a NumPy array instead of a list. The following script produces a plot identical to the previous one:

```
from matplotlib import pyplot as plt
import numpy as np

xs = np.arange(1, 6)
tops = np.arange(2, 12, 2)

plt.bar(xs, tops)
plt.show()
```

```
plt.bar(xs, tops)
plt.show()
```

The `bar()` function is more flexible than it lets on. For example, the first argument doesn't need to be a list of numbers. It could be a list of strings representing categories of data.

Suppose you wanted to plot a bar chart representing the data contained in the following dictionary:

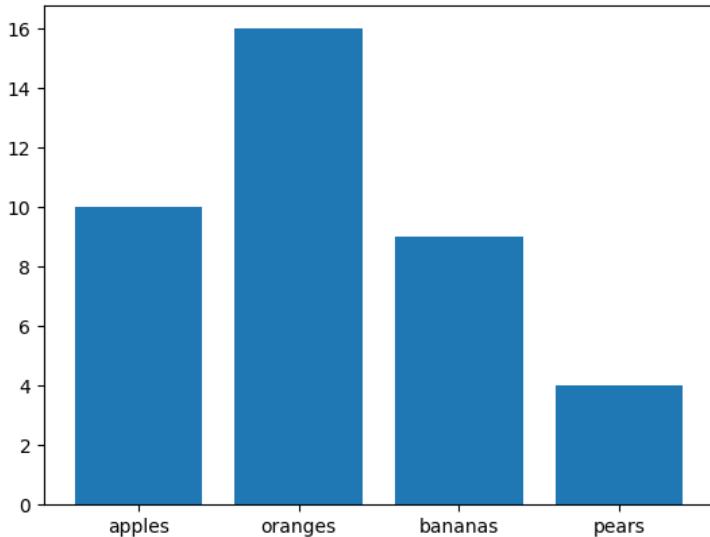
```
fruits = {
    "apples": 10,
    "oranges": 16,
    "bananas": 9,
    "pears": 4,
}
```

You can get a list of the names of the fruits using `fruits.keys()`, and the corresponding values using `fruits.values()`. Check out what happens when you pass these to the `bar()` function

```
from matplotlib import pyplot as plt

fruits = {
    "apples": 10,
    "oranges": 16,
    "bananas": 9,
    "pears": 4,
}

plt.bar(fruits.keys(), fruits.values())
plt.show()
```



The names of the fruits are conveniently used as the tick labels along the x-axis.

Note

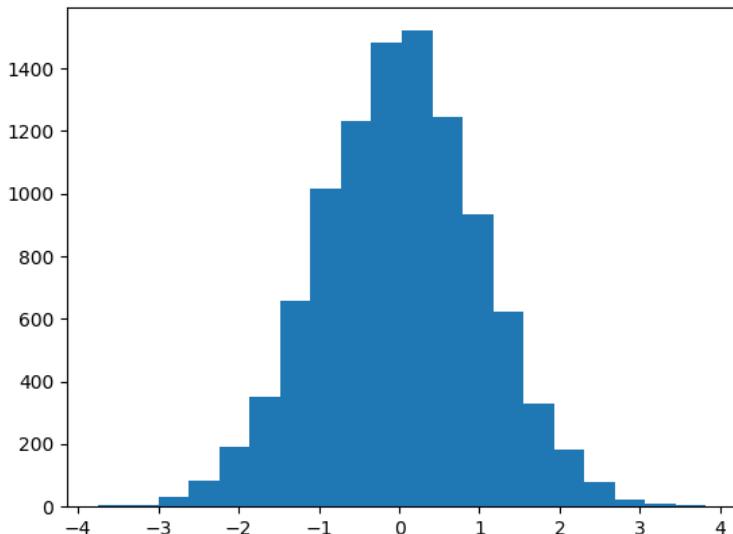
Using a list of strings as x-values works for the `plot()` function as well, although it often makes less sense to do so.

Another commonly used type of graph is the [histogram](#), which shows how data is distributed. You can make simple histograms easily with the `plt.hist()` function. You must supply `hist()` with a list (or array) of values and a number of bins to use.

For instance, we can create a histogram of 10,000 normally distributed random numbers binned across 20 possible bars with the following, which uses NumPy's `random.randn()` function to generate an array of normally distributed random numbers:

```
from matplotlib import pyplot as plt
from numpy import random

plt.hist(random.randn(10000), 20)
plt.show()
```



Note

For a detailed discussion of creating histograms with Python, check out [Python Histogram Plotting: NumPy, Matplotlib, Pandas & Seaborn on Real Python](#).

Save Figures as Images

You may have noticed that the window displaying your plots has a toolbar at the bottom. You can use this toolbar to save your plot as an image file.

More often than not, you probably don't want to have to sit at your computer and click on the save button for each plot you want to export. Fortunately, `matplotlib` makes it easy to save your plots programmatically.

To save your plot, use the `plt.savefig()` function. Pass the path to where you would like to save your plot as a string. The example below saves a simple bar chart as `bar.png` to the current working directory. If you would like to save to somewhere else, you must provide an absolute path.

```
from matplotlib import pyplot as plt
import numpy as np

xs = np.arange(1, 6)
tops = np.arange(2, 12, 2)

plt.bar(xs, tops)
plt.savefig("bar.png")
```

Note

If you want to both save a figure and display it on the screen, make sure that you save it first before displaying it!

The `show()` function pauses execution of your code and closing the display window destroys the graph, so trying to save the figure after calling `show()` results in an empty file.

Work With Plots Interactively

When you are initially tweaking the layout and formatting of a particular graph, it can be helpful to change parts of the graph without having to re-run an entire script just to see the results.

One of the easiest ways to do this is with a [Jupyter Notebook](#), which creates an interactive Python interpreter session that runs in your

browser.

Jupyter notebooks have become a staple for interacting with and exploring data, and work great with both NumPy and `matplotlib`.

For an interactive tutorial on how to use Jupyter Notebooks, check out Jupyter's [IPython In Depth](#) tutorial.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Recreate as many of the graphs shown in this section as you can by writing your own scripts without referring to the provided code.
2. It is a [well-documented fact](#) that the number of pirates in the world is correlated with a rise in global temperatures. Write a script `pirates.py` that visually examines this relationship:
 - Read in the file `pirates.csv` from the Chapter 17 practice files folder.
 - Create a line graph of the average world temperature in degrees Celsius as a function of the number of pirates in the world—that is, graph Pirates along the x-axis and Temperature along the y-axis.
 - Add a graph title and label your graph's axes.
 - Save the resulting graph out as a PNG image file.
 - Bonus: Label each point on the graph with the appropriate Year. You should do this programmatically by looping through the actual data points rather than specifying the individual position of each annotation.

[Leave feedback on this section »](#)

17.3 Summary and Additional Resources

In this chapter, you learned about two packages commonly used in the Python scientific computing stack.

In the first section, “Use NumPy for Matrix Manipulation,” you learned about the NumPy package. NumPy is used for working with multi-dimensional arrays of data. It introduces the `ndarray` object, which is commonly created using the `array` alias.

A NumPy array is a homogenous data type, meaning it can only store a single type of data. For example, a NumPy array can contain all integers, or all floats, but cannot contain both integers *and* floats. You also saw some useful functions and methods for manipulating NumPy array objects. Finally, you were introduced to the NumPy `arange()` function, which works a lot like Python’s very own `range()` function, except that returns a one-dimensions NumPy array object.

In the second section, “Use `matplotlib` for Plotting Graphs,” you learned how to use the `matplotlib` package to create simple plots using the `pyplot` interface. You built line charts, bar charts and histograms from pure Python lists and NumPy arrays using the `plot()`, `bar()` and `hist()` functions. You learned how to style and layout your plots by adding plot and axis titles, tick markers and legends.

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-16

Additional Resources

With the knowledge you gained in this chapter you should be able to work with basic data arrays and produce some simple plots. If your goal is to use Python for data science or scientific computing, you now have some foundational knowledge. To further your study, you may want to check out the following resources:

- [Real Python Data Science Tutorials](#)
- Recommended resources on [realpython.com](#)

[Leave feedback on this section »](#)

Chapter 18

Graphical User Interfaces

Throughout this book, you have been creating **command-line applications**, which are programs that are started from and produce output in a terminal window.

Command-line apps are fine for making tools that you or other developers might use, but the vast majority of software users never want to open a terminal!

Graphical User Interfaces, called **GUIs** for short and pronounced “gooey”, have windows with components like buttons and text fields. They provide users with a familiar and visual way to interact with a program.

In this chapter, you'll learn how to:

- Add a simple GUI to a command line application with EasyGUI
- Create full-featured GUI applications with Tkinter

Let's get started!

[Leave feedback on this section »](#)

18.1 Add GUI Elements With EasyGUI

You can use the EasyGUI library to quickly add a graphical user interface to your program. EasyGUI is somewhat limited, but works well for simple tools that just needs a little bit of input from the user.

In this section, you'll use EasyGUI to create a short GUI program that allows a user to pick a PDF file from their hard drive and rotate its pages by a selected amount.

Installing EasyGUI

To get started, you need to install EasyGUI with `pip3`:

```
$ pip3 install easygui
```

Once EasyGUI is installed, you can check out some details of the package with `pip3 show`:

```
$ pip3 show easygui
Name: easygui
Version: 0.98.1
Summary: EasyGUI is a module for very simple, very easy GUI
        programming in Python. EasyGUI is different from other
        GUI generators in that EasyGUI is NOT event-driven.
        Instead, all GUI interactions are invoked by simple
        function calls.

Home-page: https://github.com/robertlugg/easygui
Author: easygui developers and Stephen Ferg
Author-email: robert.lugg@gmail.com
License: BSD
Location: c:\realpython\venv\lib\site-packages
Requires:
Required-by:
```

The code in this chapter is written using EasyGUI version 0.98.1, the same version you see in the information shown above.

Your First EasyGUI Application

EasyGUI is great for displaying **dialog boxes** to collect user input and display output. It is not particularly great for creating a large application with several windows, menus, and toolbars.

You can think of EasyGUI as a sort of replacement for the `input()` and `print()` functions that you have been using for input and output.

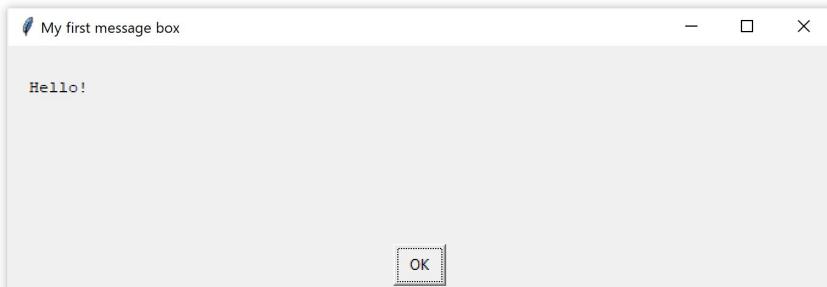
Program flow with EasyGUI typically works like this:

1. At some point in the code, a visual element is displayed on the user's screen.
2. Execution of the code pauses until the user provides input with the visual element.
3. The user's input is returned as an object and execution of the code is resumed.

To get a feel for how EasyGUI works, open a new interactive window in IDLE and execute the following lines of code:

```
>>> import easygui as gui  
>>> gui.msgbox(msg="Hello!", title="My first message box")
```

If you run the code on Windows, you will see a window like the following displayed on your screen:



The window's appearance depends on the operating system on which the code is executed. On macOS, the window looks like this one:



Here's what the window looks like on Ubuntu:



For the rest of this section, Windows screenshots will be shown.

Both EasyGUI and IDLE are written using the Tkinter library, which you'll learn about in the next section. This overlap sometimes causes issues with execution, such as dialog boxes getting frozen or stuck.

If you think this might be happening to you, try running your code from a terminal. You can start an interactive Python session from a terminal with the `python` command on Windows and `python3` on macOS/Ubuntu.

Let's break down what you see in the dialog box you generated with the code above:

1. The string "Hello!" passed to the `msg` parameter of `msgbox()` is displayed as the message in the message box.
2. The string "My first message box" passed to the `title` parameter is displayed as the title of the message box.
3. There is one button in the message box labelled `OK`.

Press the `OK` button to close the dialog box and look at IDLE's interactive window. The string '`OK`' is displayed below the last line of code you typed:

```
>>> gui.msgbox(msg="Hello, EasyGUI!", title="My first message box")
'OK'
```

`msgbox()` returns the button label when the dialog box is closed. If the dialog box is closed without pressing the `OK` button, then the value `None` is returned.

You can customize the button label by setting a third optional parameter called `ok_button`. For example, the following creates a message box with a button labeled *Click me*:

```
>>> gui.msgbox(msg="Hello!", title="Greeting", ok_button="Click me")
```

`msgbox()` is great for displaying a message, but it doesn't provide the user with many options for interacting with your program. EasyGUI has several functions that display various types of dialog boxes. Let's explore some of these now!

EasyGUI's Ensemble of GUI Elements

Besides `msgbox()`, EasyGUI has several other functions for displaying different kinds of dialog boxes. The following table summarizes some of the available functions:

Function	Description
<code>msgbox()</code>	A dialog box for displaying a message with a single button. It returns the label of the button.
<code>buttonbox()</code>	A dialog box with several buttons. It returns the label of the selected button.
<code>indexbox()</code>	A dialog box with several buttons. It returns the index of the selected button.
<code>enterbox()</code>	A dialog box with a text entry box. It returns the text entered.
<code>fileopenbox()</code>	A dialog box for selecting a file to be opened. It returns the absolute path to the selected file.
<code>diropenbox()</code>	A dialog box for selecting a directory to be opened. It returns the absolute path to the selected directory.
<code>filesavebox()</code>	A dialog box for saving a file. It returns the absolute path to the location for saving the file.

Let's look at each one of these individually.

buttonbox()

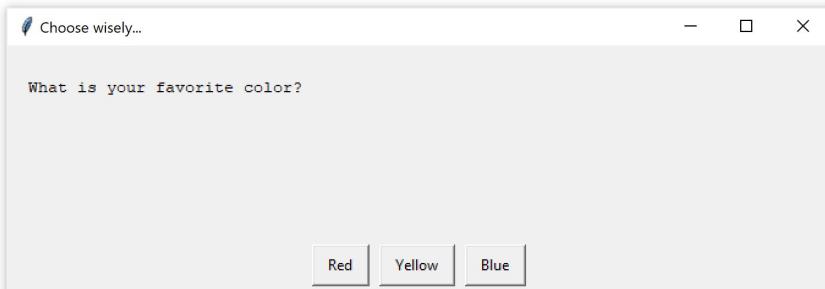
EasyGUI's `buttonbox()` displays a dialog box with a message and several buttons that the user can click. The label of the clicked button is returned to your program.

Just like `msgbox()`, the `buttonbox()` function has `msg` and `title` parameters for setting the message to be displayed and the title of the dialog box. `buttonbox()` has a third parameter called `choices` that is used to set up the buttons.

For example, the following code produces a dialog box with three buttons labelled "Red", "Yellow", and "Blue":

```
>>> gui.buttonbox(  
...     msg="What is your favorite color?",  
...     title="Choose wisely...",  
...     choices=("Red", "Yellow", "Blue"),  
... )
```

Here's what the dialog box looks like:



When you press one of the buttons, the button label is returned as a string. For example, if you press the **Yellow** button, you'll see the string 'Yellow' displayed in the output of the interactive window just below the `buttonbox()` function:

```
>>> gui.buttonbox(  
...     msg="What is your favorite color?",  
...     title="Choose wisely...",  
...     choices=("Red", "Yellow", "Blue"),  
... )  
'Yellow'
```

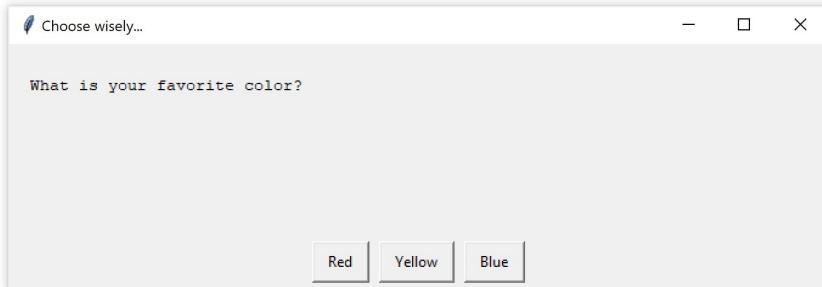
Just like `msgbox()`, the value `None` is returned if the dialog box is closed without pressing one of the buttons.

indexbox()

`indexbox()` displays a dialog box that looks identical to the dialog box displayed by `buttonbox()`. In fact, you create an `indexbox()` the same way as you do a `buttonbox()`:

```
>>> gui.indexbox(  
...     msg="What's your favorite color?",  
...     title="Choose wisely...",  
...     choices=("Red", "Yellow", "Blue"),  
... )
```

Here's what the dialog box looks like:



The difference between `indexbox()` and `buttonbox()` is that `indexbox()` returns the index of the button label in the list or tuple passed to `choices`, instead of the label itself.

For example, if you click on the `Yellow` button, the integer `1` is returned:

```
>>> gui.indexbox(  
...     msg="What's your favorite color?",  
...     title="Favorite color",  
...     choices=("Red", "Yellow", "Blue"),  
... )
```

1

Because `indexbox()` returns an index and not a string, it is a good idea to define the tuple for `choices` outside of the function so that you can reference the label by index later in your code:

```
>>> colors = ("Red", "Yellow", "Blue")
>>> choice = gui.indexbox(
    msg="What's your favorite color?",
    title="Favorite color",
    choices=colors,
)
>>> choice
1
>>> colors[choice]
'Yellow'
```

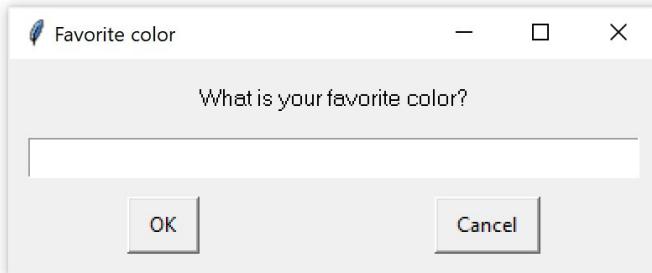
`buttonbox()` and `indexbox()` are great for getting input from a user when they need to choose from a pre-determined set of choices. These functions are not well suited to getting information such as a user's name or email address. For that, you can use the `enterbox()`.

enterbox()

`enterbox()` is used to collect text input from a user:

```
>>> gui.enterbox(
...     msg="What is your favorite color?",
...     title="Favorite color",
... )
```

The dialog box produced by `enterbox()` has an input box where the user can type in their own answer:



Type in a color name, such as `Yellow`, and press `OK`. The text you entered is returned as a string:

```
>>> gui.enterbox(  
...     msg="What is your favorite color?",  
...     title="Favorite color",  
... )  
'Yellow'
```

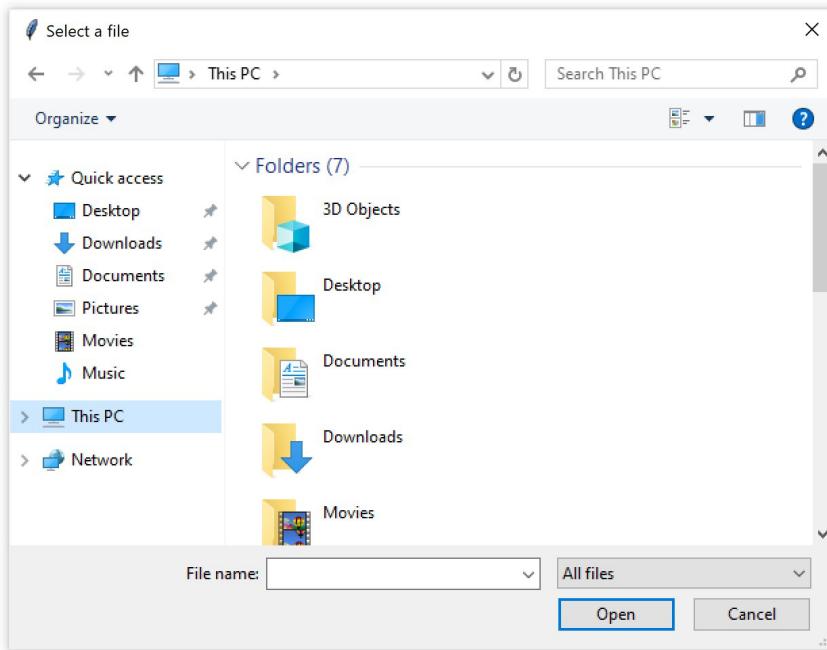
One of the most common reasons for displaying a dialog box is to allow a user to select a file or folder in their filesystem. EasyGUI has some special functions designed just for these operations.

fileopenbox()

`fileopenbox()` displays a dialog box for selecting a file to be opened:

```
>>> gui.fileopenbox(title="Select a file")
```

The dialog box looks like the standard system file open dialog box:



Select a file and click the `Open` button. A string containing the full path to the selected file is returned.

Important

`fileopenbox()` does not actually open the file! To do that you need to use `theopen()` built-in like you learned to do in Chapter 12.

Just like `msgbox()` and `buttonbox()`, the value `None` is returned if the user presses `Cancel` or closes the dialog box without selecting a file.

diropenbox() and filesavebox()

EasyGUI has two other functions that generate dialogs nearly identical to the one generated by `fileopenbox()`:

1. `diropenbox()` opens a dialog that can be used to select a folder in-

stead of a file. When the user presses **Open**, the full path to the directory is returned.

2. `filesavebox()` opens a dialog to select a location for saving a file and will confirm that the user wants to overwrite the file if the chosen name already exists. Just like `fileopenbox()`, the file path is returned when the user presses **Save**. The file is not actually saved.

Important

Neither `diropenbox()` and `filesavebox()` actually open a directory or save a file. They only return the absolute path to the directory to opened or the file to be saved.

You must write the code yourself to open the directory or save the file.

Both `diropenbox()` and `filesavebox()` return `None` if the dialogs are closed without pressing **Open** or **Save**. This can cause your program to crash if you aren't careful.

For example, the following raises a `TypeError` if the dialog box is closed without making any selection:

```
>>> path = gui.fileopenbox(title="Select a file")
>>> open_file = open(path, "r")
Traceback (most recent call last):
  File "<stdin>", line 2, in <module>
TypeError: expected str, bytes or os.PathLike object, not NoneType
```

How you handle situations like these has a huge impact on a user's experience with your program.

Exiting Your Program Gracefully

Suppose you are writing a program for extracting pages from a PDF file. The first thing the program might do is use `fileopenbox()` so that the user can select with PDF to open.

What do you do if the user decides they don't want to run the program and presses the **Cancel**?

You must make sure that your program handles these situations gracefully. The program shouldn't crash or produce any unexpected output. In the situation described above, the program should stop probably just stop running altogether.

One way to stop a program from running is with Python's built-in `exit()` function.

For example, the following program uses `exit()` to stop the program when the user presses the **Cancel** button in a file selection dialog box:

```
import easygui as gui

path = gui.fileopenbox(title="Select a file")

if path is None:
    exit()
```

If the user closes the file open dialog box without pressing **OK**, then `path` is `None` and the program executes the `exit()` function in the `if` block. This program closes and execution stops.

Note

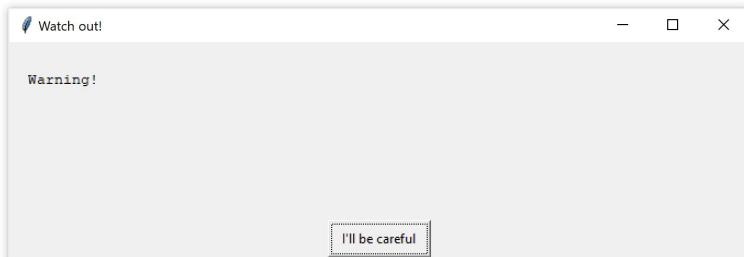
If you're running the program in IDLE, `exit()` also closes the current interactive window. It's very thorough.

Now that you know how to create dialog boxes with EasyGUI, let's put everything together into a real-world application.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Create the following dialog box:



2. Create the following dialog box:



[Leave feedback on this section »](#)

18.2 Example App: PDF Page Rotator

EasyGUI is a great choice for utility applications that automate simple yet repetitive tasks. If you work in an office, you can really boost your productivity by creating tools with EasyGUI that take the pain out of everyday TODO items.

In this section, you'll use some of the EasyGUI dialog boxes you learned about in the last section to create an application for rotating PDF pages.

In doing so, you'll bring together a lot of the concepts you've learned about thus far, including for loops (Chapter 6), conditional logic

(Chapter 8), reading and writing files (Chapter 12), and working with PDF files (Chapter 14).

The Application Design

Before we dive into the code, let's put some thought into how the program should work.

The program needs to ask the user which PDF file to open, by how many degrees they want to rotate each page, and where the user would like to save the new PDF. Then the program needs to open the file, rotate the pages, and save the new file.

Let's map this out into explicit steps that we can more easily translate into code:

1. Display a file selection dialog for opening a PDF file.
2. If the user cancels the dialog, then exit the program.
3. Let the user select one of 90, 180 or 270 degrees to rotate the PDF pages.
4. Display a file selection dialog for saving the rotated PDF.
5. If the user tries to save a file with the same name as the input file:
 - Alert the user with a message box that this is not allowed.
 - Return to step 4.
6. If the user cancels the file save dialog, then exit the program.
7. Perform the page rotation:
 - Open the selected PDF.
 - Rotate all of the pages.
 - Save the rotated PDF to the selected file.

Note

When you are designing an application, it helps to plan out each step before you start coding. For large applications, drawing diagrams describing the program flow can help keep everything organized.

Implementing the Design

Now that we have a plan, let's tackle each step one at a time. Open a new script window in IDLE to follow along.

First, import EasyGUI and PyPDF2:

```
import easygui as gui  
from PyPDF2 import PdfFileReader, PdfFileWriter
```

Step 1 in our plan is to display a file selection dialog for opening a PDF file. We can do this with `fileopenbox()`:

```
# 1. Display a file selection dialog for opening a PDF file.  
input_path = gui.fileopenbox(  
    title="Select a PDF to rotate...",  
    default="*.pdf"  
)
```

Here we've set the `default` parameter to `"*.pdf"`, which configures the dialog to only display files with the `.pdf` extension. This helps prevent the user from accidentally selecting a file that isn't a PDF.

The file path selected by the user is assigned to the `input_path` variable. If the user closed the dialog without selecting a file path (Step 2), then `input_path` is `None`. In this case, we need to exit the program:

```
# 2. If the user cancels the dialog, then exit the program.  
if input_path is None:  
    exit()
```

The third step is to ask the user how much they would like to rotate the PDF pages. They can choose either 90, 180, or 270 degrees. Let's use a `buttonbox()` to collect this information:

```
# 3. Let the user select one of '90', '180' or '270' degrees to rotate
# the PDF pages.

choices = ("90", "180", "270")
degrees = gui.buttonbox(
    msg="Rotate the PDF clockwise by how many degrees?",
    title="Choose rotation...",
    choices=choices,
)
)
```

The dialog generated here has three buttons with the labels "90", "180", and "270". When the user clicks on one of these buttons, the label of the button is assigned to the `degrees` variable as a string.

In order to rotate the pages in the PDF by the selected angle, we'll need the value to be an integer, not a string. Let's go ahead and convert it to an integer:

```
degrees = int(degrees)
```

Next, get the output file path from the user using `filesavebox()`:

```
# 4. Display a file selection dialog for saving the rotated PDF.

save_title = "Save the rotated PDF as..."
file_type = "*.pdf"
output_path = gui.filesavebox(title=save_title, default=file_type)
```

Just like `fileopenbox()`, we've set the `default` parameter to `*.pdf`. This ensures that the file automatically gets saved with the `.pdf` extension.

The user shouldn't be allowed to overwrite the original file (Step 5). You can use a `while` loop to repeatedly show the user a warning until they pick a path that is different from the input file path:

```
# 5. If the user tries to save with the same name as the input file:  
while input_path == output_path:  
    # - Alert the user with a message box that this is not allowed.  
    gui.msgbox(msg="Cannot overwrite original file!")  
    # - Return to step 4.  
    output_path = gui.filesavebox(title=save_title, default=file_type)
```

The while loop checks if `input_path` is the same as `output_path`. If it isn't, then the loop body is ignored. If `input_path` and `output_path` are the same, then `msgbox()` is used to show a warning to the user telling them they can't overwrite the original file.

After warning the user, `filesavebox()` is used to display another file save dialog box with the same title and default file type as before. This is the part that returns the user to step 4. Even though the program doesn't actually return the line of code where `filesavebox()` is first called, the effect is the same.

If the user closes the file save dialog without pressing `Save`, the program should exit (Step 6):

```
# 6. If the user cancels the file save dialog, then exit the program.  
if output_path is None:  
    exit()
```

Now you have everything you need to implement the last step of the program:

```
# 7. Perform the page rotation:  
#     - Open the selected PDF.  
input_file = PdfFileReader(input_path)  
output_pdf = PdfFileWriter()  
  
#     - Rotate all of the pages.  
for page in input_file.pages:  
    page = page.rotateClockwise(degrees)  
    output_pdf.addPage(page)
```

```
#     - Save the rotated PDF to the selected file.  
with open(output_path, "wb") as output_file:  
    output_pdf.write(output_file)
```

Try out your new PDF rotation application! It works equally well on Windows, macOS, and Ubuntu Linux!

Here's the full application source code for your reference:

```
import easygui as gui  
from PyPDF2 import PdfFileReader, PdfFileWriter  
  
# 1. Display a file selection dialog for opening a PDF file.  
input_path = gui.fileopenbox(  
    title="Select a PDF to rotate...",  
    default="*.pdf"  
)  
  
# 2. If the user cancels the dialog, then exit the program.  
if input_path is None:  
    exit()  
  
# 3. Let the user select one of `90`, `180` or `270` degrees to rotate  
# the PDF pages.  
choices = ("90", "180", "270")  
degrees = gui.buttonbox(  
    msg="Rotate the PDF clockwise by how many degrees?",  
    title="Choose rotation...",  
    choices=choices,  
)  
  
# 4. Display a file selection dialog for saving the rotated PDF.  
save_title = "Save the rotated PDF as..."  
file_type = "*.pdf"  
output_path = gui.filesavebox(title=save_title, default=file_type)
```

```
# 5. If the user tries to save with the same name as the input file:  
while input_path == output_path:  
    # - Alert the user with a message box that this is not allowed.  
    gui.msgbox(msg="Cannot overwrite original file!")  
    # - Return to step 4.  
    output_path = gui.filesavebox(title=save_title, default=file_type)  
  
# 6. If the user cancels the file save dialog, then exit the program.  
if output_path is None:  
    exit()  
  
# 7. Perform the page rotation:  
#     - Open the selected PDF.  
input_file = PdfFileReader(input_path)  
output_pdf = PdfFileWriter()  
  
#     - Rotate all of the pages.  
for page in input_file.pages:  
    page = page.rotateClockwise(degrees)  
    output_pdf.addPage(page)  
  
#     - Save the rotated PDF to the selected file.  
with open(output_path, "wb") as output_file:  
    output_pdf.write(output_file)
```

EasyGUI is great for quickly creating a GUI for small tools and applications. For larger projects, EasyGUI may be too limited. That's where Python's built-in `Tkinter` library comes in.

Tkinter is a GUI framework that operates at a lower level than EasyGUI. That means you have more control over the visual aspects of the GUI, such as window size, font size, font color, and what GUI elements are present in a dialog box or window.

The rest of this chapter is devoted to developing GUI applications with Python's built-in Tkinter library.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. The GUI application for rotating PDF pages in this section has a problem. The program crashes if the user closes the `buttonbox()` used to select degrees without selecting a value.

Fix this problem by using a `while` loop to keep displaying the selection dialog if `degrees` is `None`.

[Leave feedback on this section »](#)

18.3 Challenge: PDF Page Extraction Application

In this challenge, you'll use EasyGUI to write a GUI application for extracting pages from a PDF file.

Here's a detailed plan for the application:

1. Ask the user to select a PDF file to open.
2. If no PDF file is chosen, exit the program.
3. Ask for a starting page number.
4. If the user does not enter a starting page number, exit the program.
5. Valid page numbers are positive integers. If the user enters an invalid page number:
 - Warn the user that the entry is invalid .
 - Return to step 3.
6. Ask for an ending page number.
7. If the user does not enter an ending page number, exit the program.
8. If the user enters an invalid page number:

- Warn the user that the entry is invalid .
 - Return to step 6.
9. Ask for the location to save the extracted pages.
 10. If the user does not select a save location, exit the program.
 11. If the chosen save location is the same as the input file path:
 - Warn the user that they can not overwrite the input file.
 - Return to step 9.
 12. Perform the page extraction:
 - Open the input PDF file.
 - Write a new PDF file containing only the pages in the selected page range.

You can find the solutions to this code challenge and many other bonus resources online at realpython.com/python-basics/resources.

[Leave feedback on this section »](#)

18.4 Introduction to Tkinter

Python has [a lot](#) of GUI frameworks, but [Tkinter](#) is the only framework that is built into the Python standard library.

Tkinter has several strengths. It is **cross-platform**, meaning the same code works on Windows, macOS, and Linux. Visual elements are rendered using native operating system elements, so applications built with Tkinter look like they belong on the platform where they are run.

Although Tkinter is considered the *de facto* Python GUI framework, it is not without criticism. One notable criticism is that GUIs built with Tkinter look outdated. If you want a shiny, modern interface, then Tkinter may not be what you are looking for.

However, Tkinter is lightweight and is relatively simple to use compared to other frameworks. This makes it a compelling choice for

building GUI applications in Python, especially for applications where a modern sheen is unnecessary and quickly building something that is functional and cross-platform is the top priority.

Note

As was mentioned in the last section, IDLE is built with Tkinter. You may encounter difficulties when running your own GUI programs within IDLE.

If you find that the GUI window you are trying to create is unexpectedly freezing or appears to be making IDLE misbehave in some unexpected way, try running your script from the a command prompt or terminal.

Let's dive right in and see how you build an application with Tkinter.

Your First Tkinter Application

The foundational element of a Tkinter GUI is the **window**. Windows are the containers in which all other GUI elements live. Other GUI elements, such as text boxes, labels, and buttons, are known as **widgets**. Widgets are contained inside of windows.

Let's create a window that contains a single widget. Start by opening a new interactive window in IDLE.

The first thing you need to do is import the Tkinter module:

```
>>> import tkinter as tk
```

A window is an instance of Tkinter's `Tk` class. Go ahead and create a new window and assign it to the variable `window`:

```
>>> window = tk.Tk()
```

When you execute the above code, a new window pops up on your screen. How it looks depends on your operating system:



For the rest of this chapter, Windows screenshots will be used.

Now that we have a window, let's add a widget. The `tk.Label` class is used to add some text to a window.

Create a `Label` widget with the text "Hello, Tkinter" and assign it to a variable called `greeting`:

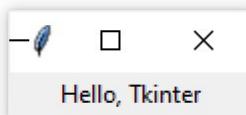
```
>>> greeting = tk.Label(text="Hello, Tkinter")
```

The window you created earlier doesn't change. You just created a `Label` widget, but it hasn't been added to the window yet.

There are several ways to add widgets to a window. Right now, we'll use the `Label` widget's `.pack()` method:

```
>>> greeting.pack()
```

The window now looks like this:



When you `.pack()` a widget into a window, Tkinter sizes the window as small as it can while still fully encompassing the widget.

Now execute the following:

```
>>> window.mainloop()
```

Nothing seems to happen, but notice that a new prompt does not appear in the shell.

`window.mainloop()` tells Python to run the Tkinter application and **blocks** any code that comes after it from running until the window it's called on is closed. Go ahead and close the window you've created and you'll see a new prompt displayed in the shell.

Important

When you work with Tkinter from a REPL like IDLE's interactive window, updates to windows are applied as each line is executed.

This is not the case when a Tkinter program is executed from a Python file.

If you do not include `window.mainloop()` at the end of a program in a Python file, the Tkinter application will never run, and nothing will be displayed.

Creating a window with Tkinter only takes a couple of lines of code. But blank windows aren't very useful! In the next section, you'll learn about some of the widgets available in Tkinter, and how you can customize them to meet your application's needs.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Using Tkinter from IDLE's interactive window, execute code that creates a window with a `Label` widget with the text "GUIs are great!".
2. Repeat Exercise 1 with the text "Python rocks!".

3. Repeat Exercise 1 with the text "Engage!".

Leave feedback on this section »

18.5 Working With Widgets

Widgets are the bread and butter of Tkinter. They are the elements through which users interact with your program.

Each widget in Tkinter is defined by a class. Here are some of the widgets available:

Widget Class	Description
Label	A widget used to display text on the screen.
Button	A button that can contain text and can perform an action when clicked.
Entry	A text entry widget that allows only a single line of text.
Text	A text entry widget that allows multiline text entry.
Frame	A rectangular region used to group related widgets or provide padding between widgets.

You'll see how to work with each of these in the following sections.

Note

Tkinter has many more widgets than the ones listed here. For a full list, check out the [Basic Widgets](#) and [More Widgets](#) articles in the [TkDocs](#) tutorial.

Let's take a closer look at the `Label` widget.

Label Widgets

`Label` widgets are used to display text or images. The text displayed by a `Label` widget can't be edited by the user. It is for display purposes

only.

As you saw in the example at the beginning of this chapter, you can create a `Label` widget by instantiating the `Label` class and passing a string to the `text` parameter:

```
label = tk.Label(text="Hello, Tkinter")
```

`Label` widgets display text with the default system text color and the default system text background color. These are typically black and white, respectively, but you may see different colors if you have changes these settings in your operating system.

You can control `Label` text and background colors using the `foreground` and `background` parameters:

```
label = tk.Label(  
    text="Hello, Tkinter",  
    foreground="white", # Set the text color to white  
    background="black" # Set the background color to black  
)
```

There are numerous valid color name, including:

- "red"
- "orange"
- "yellow"
- "green"
- "blue"
- "purple"

Many of the [HTML color names](#) work with Tkinter.

Note

A chart with most of the valid color names is available [here](#). For a full reference, including macOS and Windows-specific system colors that are controlled by the current system theme, check out [this list](#).

You can also specify a color using hexadecimal RGB values:

```
label = tk.Label(text="Hello, Tkinter", background="#34A2FE")
```

This sets the label background to a nice light blue color.

Hexadecimal RGB values are more cryptic than named colors, but they are more flexible. Fortunately, there are [tools](#) available that make getting hexadecimal color codes relatively painless.

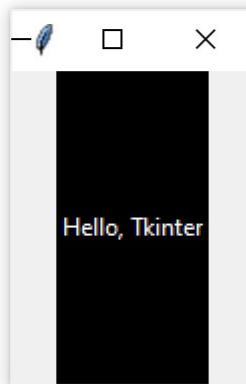
If you don't feel like typing out `foreground` and `background` all the time, you can use the shorthand `fg` and `bg` parameters to set the foreground and background colors:

```
label = tk.Label(text="Hello, Tkinter", fg="white", bg="black")
```

You can also control the width and height of a label with the `width` and `height` parameters:

```
label = tk.Label(  
    text="Hello, Tkinter",  
    fg="white",  
    bg="black",  
    width=10,  
    height=10  
)
```

Here's what this label looks like in a window:



It may seem strange that the label in the window is not square even though width and height are both set to 10. This is because the height and width are measured in **text units**.

One horizontal text unit is determined by the width of the character "0" (the number zero) in the default system font. Similarly, one vertical text unit is determined by the height of the character "0".

Note

Tkinter uses text units for width and height measurements, instead of something like inches, centimeters, or pixels, to ensure consistent behavior of the application across platforms.

Measuring units by the width of a character means that the size of a widget is relative to the default font on a user's machine. This ensures text fits properly in labels and buttons, no matter where the application is running.

Labels are great for displaying some text, but they don't help you get input from a user. The next three widgets that we'll look at are all used to get user input.

Button Widgets

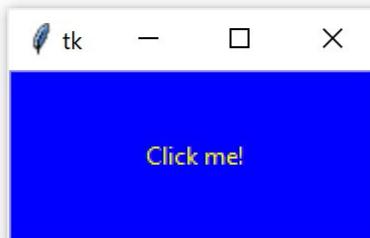
Button widgets are used to display clickable buttons. They can be configured to call a function whenever they are clicked. We'll talk about how to call functions from button clicks in the next section. For now, let's look at how to create and style a Button.

There are many similarities between `Button` and `Label` widgets. In many ways, a `Button` is just a `Label` that you can click! The same keyword arguments used to create and style a `Label` work with `Button` widgets.

For example, the following code creates a `Button` with a blue background, yellow text, and height and width set to 10 and 5 text units, respectively:

```
button = tk.Button(  
    text="Click me!",  
    width=25,  
    height=5,  
    bg="blue",  
    fg="yellow",  
)
```

Here's what the button looks like in a window:



Pretty nifty!

The next two widgets we'll see are used to collect text input from a user.

Entry Widgets

When you need to get a little bit of text from a user, like a name or an email address, use an `Entry` widget. They display a small text box that the user can type some text into.

Creating and styling an `Entry` widget works pretty much exactly like `Label` and `Button` widgets. For example, the following creates a widget with a blue background, yellow text, and a width of 50 text units:

```
entry = tk.Entry(fg="yellow", bg="blue", width=50)
```

The interesting bit about `Entry` widgets isn't how to style them, though. It's how to use them get input from a user. There are three main operations that you can perform with `Entry` widgets:

1. Retrieving text with the `.get()` method
2. Deleting text with the `.delete()` method
3. Inserting text with the `.insert()` method

The best way to get a grip on `Entry` widgets is to create one and interact with it. Go ahead and open IDLE's interactive window and follow along with the examples in this section.

First, import `tkinter` and create a new window:

```
>>> import tkinter as tk  
>>> window = tk.Tk()
```

Now create a `Label` and an `Entry` widget:

```
>>> label = tk.Label(text="Name")  
>>> entry = tk.Entry()
```

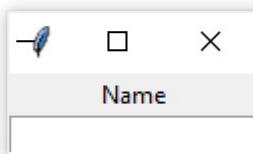
The `Label` describes what sort of text should go in the `Entry` widget. It

doesn't enforce any sort of requirements on the `Entry`, but it tells the user what our program expects them to put there.

We need to `.pack()` the widgets into the window so that they are visible:

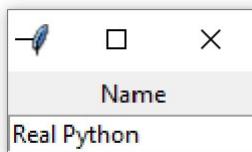
```
>>> label.pack()  
>>> entry.pack()
```

Here's what that looks like:



Notice that Tkinter automatically centers the `Label` above the `Entry` widget in the window. This is a feature of the `.pack()` method, which you'll learn more about in later sections.

Click inside the `Entry` widget with your mouse and type "Real Python":



Now you've got some text entered into the `Entry` widget, but that text hasn't been sent to your program yet.

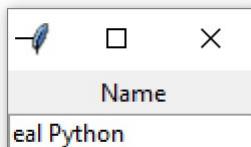
Use the `Entry` widget's `.get()` method to retrieve the text and assign it to a variable called `name`:

```
>>> name = entry.get()  
>>> name  
'Real Python'
```

You can delete text using the `Entry` widget's `.delete()` method. `.delete()` takes an integer argument that tells it which character to remove. For example, `.delete(0)` deletes the first character from the `Entry`:

```
>>> entry.delete(0)
```

The text remaining in the widget is now "eal Python":



Note

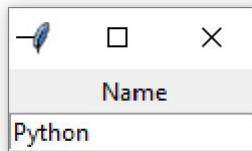
Just like Python string objects, text in an `Entry` widget is indexed starting with 0.

If you need to remove several characters from an `Entry`, pass a second integer argument to `.delete()` indicating the index of the character where deletion should stop.

For example, the following deletes the first four letters in the `Entry`:

```
>>> entry.delete(0, 4)
```

The remaining text now reads "Python":



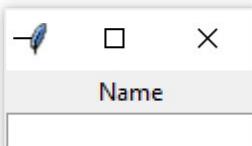
Note

`Entry.delete()` works just like string slices. The first argument determines the starting index and the deletion continues up to **but not including** the index passed as the second argument.

Use the special constant `tk.END` for the second argument of `.delete()` to remove all text in an `Entry`:

```
>>> entry.delete(0, tk.END)
```

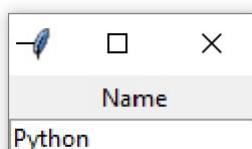
You'll now see a blank text box:



To insert text into an `Entry` widget, use the `.insert()` method:

```
>>> entry.insert(0, "Python")
```

The window now looks like this:



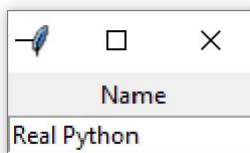
The first argument tells `.insert()` where to insert the text. If there is no text in the `Entry`, the new text will always be inserted at the beginning of the widget, no matter what value you pass to the first argument.

For example, calling `.insert()` with 100 as the first argument instead of 0, as you did above, would have generated the same output.

If an `Entry` already contains some text, `.insert()` will insert the new text at the specified position and shift all existing text to the right:

```
>>> entry.insert(0, "Real ")
```

The widget text now reads "Real Python":



`Entry` widgets are great for capturing small amounts of text from a user, but because they are only displayed on a single line, they are not ideal for gathering large amounts of text. That's where `Text` widgets come in!

Text Widgets

`Text` widgets are used for entering text, just like `Entry` widgets. The difference is that `Text` widgets may contain multiple lines of text.

With a `Text` widget, a user can input a whole paragraph, or even several pages, of text!

Just like `Entry` widgets, there are three main operations you can perform with `Text` widgets:

1. Retrieve text with the `.get()` method

2. Delete text with the `.delete()` method
3. Insert text with the `.insert()` method

Although the method names are the same as the `Entry` methods, they work a bit differently. Let's get our hands dirty by creating a `Text` widget and seeing what all it can do.

Note

If you still have the window from the previous section open, you can close it by executing the following in IDLE's interactive window:

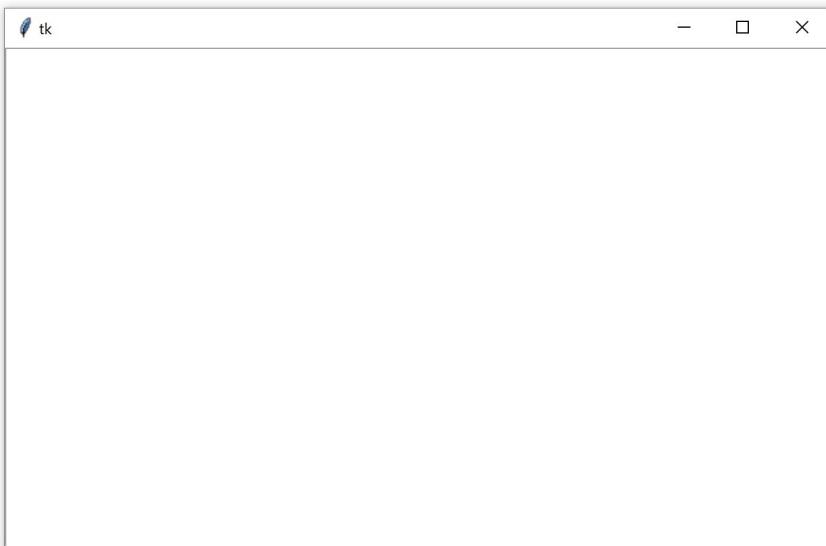
```
>>> window.destroy()
```

You can also close it manually by clicking the `Close` button on the window itself.

In IDLE's interactive window, create a new blank window and `.pack()` a `Text()` widget into it:

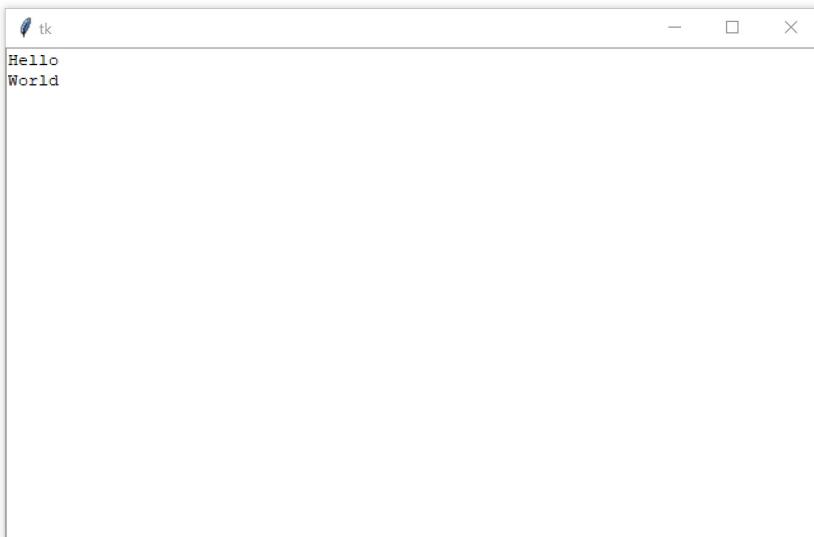
```
>>> window = tk.Tk()  
>>> text_box = tk.Text()  
>>> text_box.pack()
```

Text boxes are much larger than `Entry` widgets by default. Here's what the window created above looks like:



Click anywhere inside the window to activate the text box. Type in the word "Hello". Then press Enter and type "World" on the second line.

The window should now look like this:



Just like `Entry` widgets, you can retrieve the text from a `Text` widget using `.get()`. However, calling `.get()` with no arguments doesn't return the full text in the text box like it does for `Entry` widgets. It raises an exception:

```
>>> text_box.get()
Traceback (most recent call last):
  File "<pyshell#4>", line 1, in <module>
    text_box.get()
TypeError: get() missing 1 required positional argument: 'index1'
```

`Text.get()` required at least one argument. Calling `.get()` with a single index returns a single character. To retrieve several characters, you need to pass two arguments: a start index and an end index.

Indices in `Text` widgets work differently than `Entry` widgets. Since `Text` widgets can have several lines of text, an index must contain two pieces of information:

1. The line number of a character
2. The position of a character on that line

Line numbers start with 1 and character positions start with 0.

To make an index, you create a string of the form "<line>.<char>", replacing <line> with the line number and <char> with the character number.

For example, "1.0" represents the first character on the first line. "2.3" represents the fourth character on the second line.

Let's use the index "1.0" to get the first letter from the text box we created earlier:

```
>>> text_box.get("1.0")
'H'
```

There are five letters in the word "Hello", and the character number of o is 4, since character numbers start from 0 and the word "Hello" starts at the first position in the text box. Just like Python string slices, in order to get the entire word Hello from the text box, the end index must be one more than the index of the last character to be read.

So, to get the word "Hello" from the text box, use "1.0" for the first index and "1.5" for the second index:

```
>>> text_box.get("1.0", "1.5")
'Hello'
```

To get the word "World" on the second line of the text box, change the line numbers in each index to 2:

```
>>> text_box.get("2.0", "2.5")
'World'
```

To get all of the text in a text box, set the starting index in "1.0" and use the special tk.END constant for the second index:

```
>>> text_box.get("1.0", tk.END)
'Hello\nWorld\n'
```

Notice that text returned by .get() includes any newline characters.

You can also see from this example that every line in a `Text` widget has a newline character at the end, including the last line of text in the text box.

The `.delete()` method is used to delete characters from a text box. It works just like the `.delete()` method for `Entry` widgets.

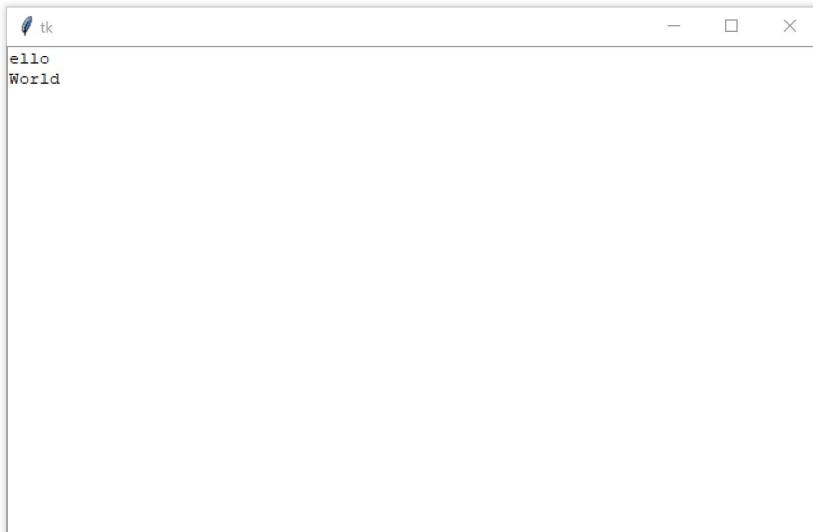
There are two ways to use the `.delete()` method:

1. With a single argument
2. With two arguments

Using the single argument version, you pass to `.delete()` the index of a single character to be deleted. For example, the following deletes the first character `h` from the text box:

```
>>> text_box.delete("1.0")
```

The first line of text in the window now reads "ello":

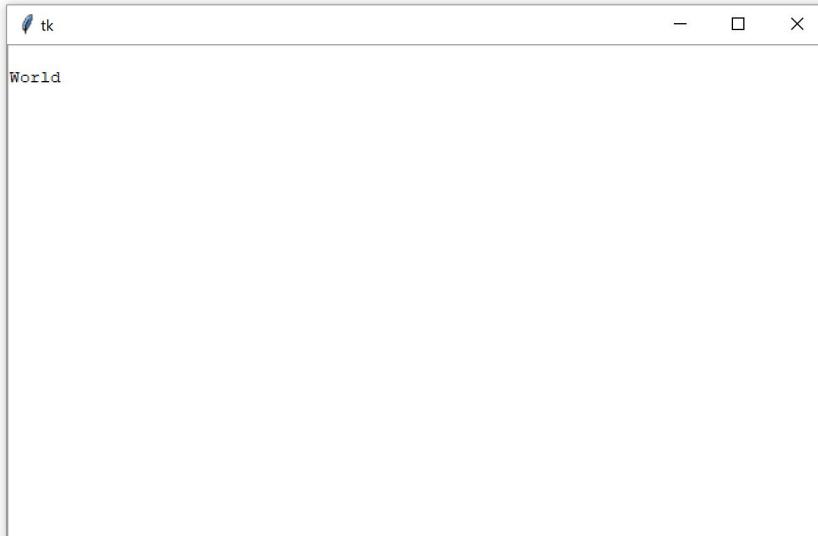


With the two argument version, you pass two indices to delete a range of characters starting at the first index and up to, but not including, the second index.

For example, to delete the remaining "ello" on the first line of the text box, use the indices "1.0" and "1.4":

```
>>> text_box.delete("1.0", "1.4")
```

Notice that the text is gone from the first line, leaving a blank line followed the word `World` on the second line:



Even though you can't see it, there is still a character on the first line. It's the newline character!

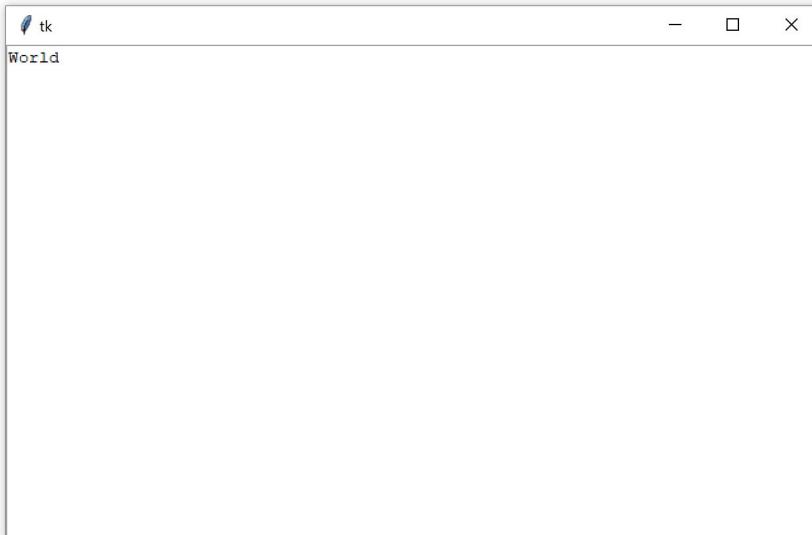
You can verify this using `.get()`:

```
>>> text_box.get("1.0")
'\n'
```

If you delete that character, the rest of the contents of the text box will shift up a line:

```
>>> text_box.delete("1.0")
```

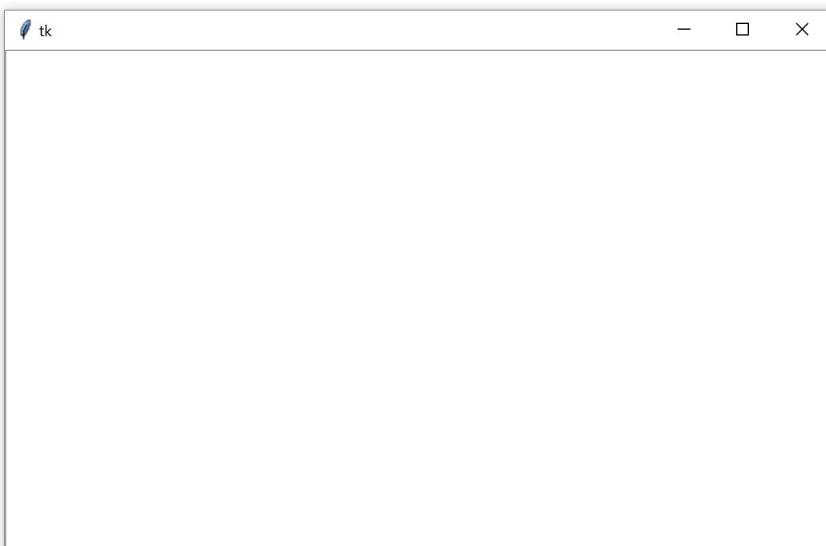
Now "World" is on the first line of the text box:



Let's clear out the rest of the text in the text box. Set "1.0" as the start index and use `tk.END` for the second index:

```
>>> text_box.delete("1.0", tk.END)
```

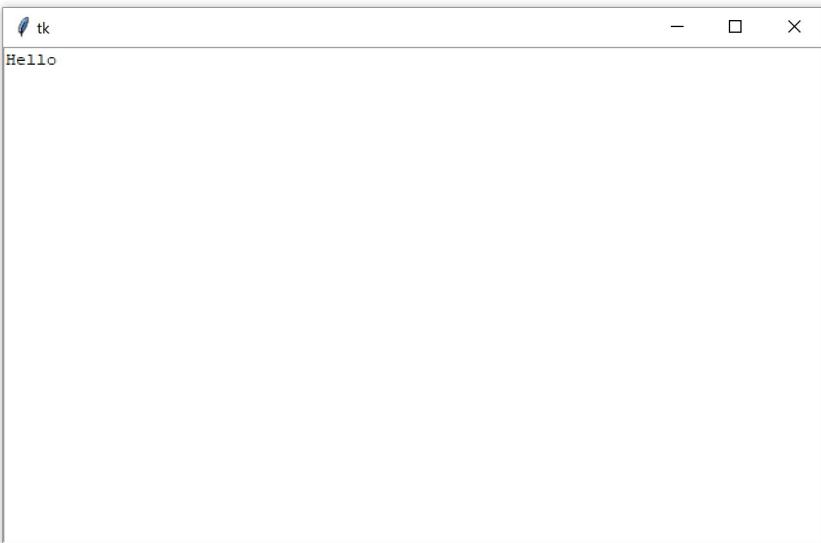
The text box is now empty:



You can insert text into a text box using the `.insert()` method:

```
>>> text_box.insert("1.0", "Hello")
```

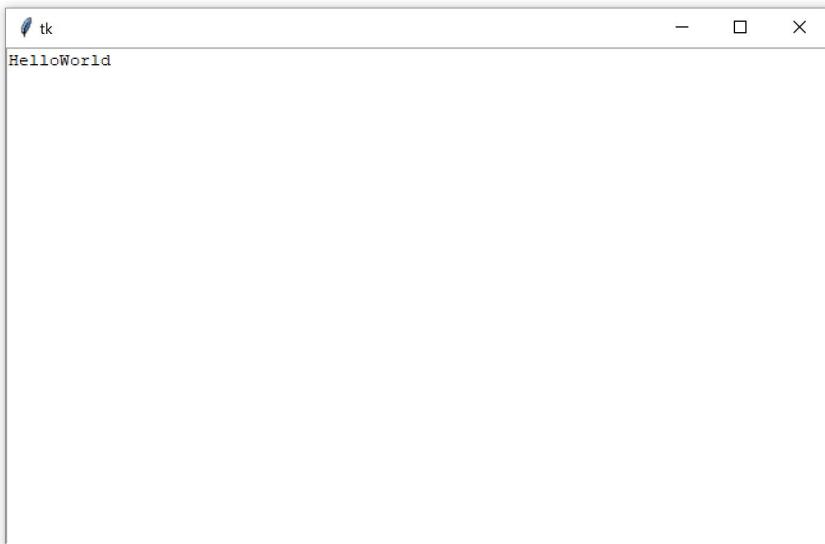
This inserts the word "Hello" at the beginning of the text box, using the same "`<line>.<column>`" format used by `.get()` to specify the insertion position:



Check out what happens if you try to insert the word "World" on the second line:

```
>>> text_box.insert("2.0", "World")
```

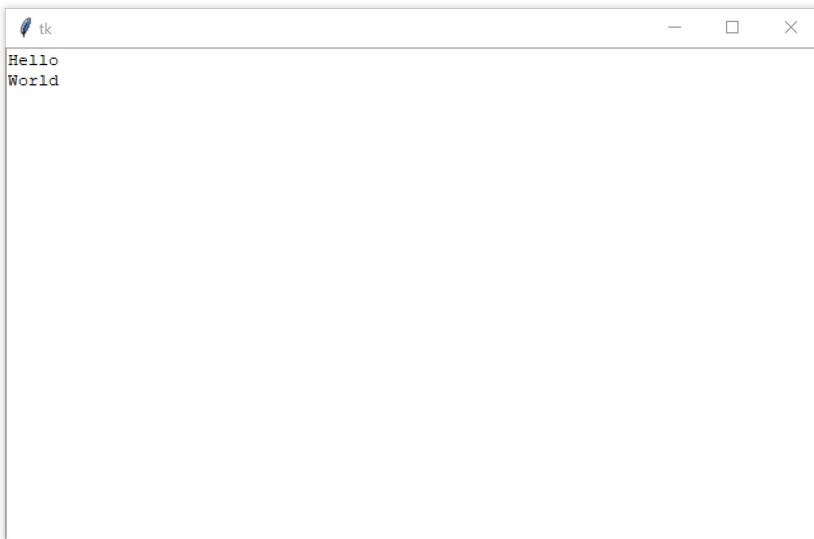
Instead of inserting the text on the second line, the text is inserted at the end of the first line:



If you want to insert text onto a new line, you need to manually insert a newline character into the string being inserted:

```
>>> text_box.insert("2.0", "\nThis goes on the second line")
```

Now "World" is on the second line of the text box:



So, `.insert()` will either insert text at the specified position, if there is already text at that position, or append text to the specified line if the character number is greater than the index of the last character in the text box.

It's usually impractical to try and keep track of what the index of the last character is. The best way to insert text at the end of a `Text` widget is pass `tk.END` to the first parameter of `.insert()`:

```
text_box.insert(tk.END, "Put me at the end!")
```

Don't forget to include the newline character `\n` at the beginning of the text if you want to put it on a new line:

```
text_box.insert(tk.END, "\nPut me on a new line!")
```

`Label`, `Button`, `Entry`, and `Text` widgets are just a few of the widgets available in Tkinter. There are several others, including widgets for checkboxes, radio buttons, scroll bars, and progress bars. For more information on the other widgets available, check out the [tutorial on tk-docs.com](#).

In this chapter, we're going to work with only five widgets: the four you have seen so far plus the `Frame` widget. `Frame` widgets are important for organizing the layout of your widgets in an application.

Before we get into the details about laying out the visual presentation of your widgets, let's take a closer look at how `Frame` widgets work, and how you can assign other widgets to them.

Assigning Widgets to Frames

The following script creates a blank `Frame` widget and assigns it to the main application window:

```
import tkinter as tk

window = tk.Tk()
frame = tk.Frame()
frame.pack()

window.mainloop()
```

The `frame.pack()` method packs the frame into the window so that the window sizes itself as small as possible to encompass the frame.

When you run the above script, you get some seriously uninteresting output:



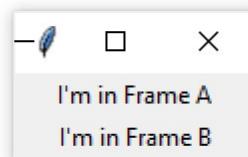
An empty `Frame` widget is practically invisible. Frames are best thought of as containers for other widgets. You can assign a widget to a frame by setting the widget's `master` attribute:

```
frame = tk.Frame()  
label = tk.Label(master=frame)
```

To get a feel for how this works, let's write a script that creates two Frame widgets called `frame_a` and `frame_b`. `frame_a` contains a label with the text "I'm in Frame A", and `frame_b` contains the label "I'm in Frame B". Here's one way to do that:

```
import tkinter as tk  
  
window = tk.Tk()  
  
frame_a = tk.Frame()  
frame_b = tk.Frame()  
  
label_a = tk.Label(master=frame_a, text="I'm in Frame A")  
label_a.pack()  
  
label_b = tk.Label(master=frame_b, text="I'm in Frame B")  
label_b.pack()  
  
frame_a.pack()  
frame_b.pack()  
  
window.mainloop()
```

Notice that `frame_a` is packed into the window before `frame_b`. The window that opens shows the label in `frame_a` above the label in `frame_b`:



Now let's see what happens when you swap the order of `frame_a.pack()`

and `frame_b.pack()`:

```
import tkinter as tk

window = tk.Tk()

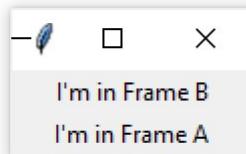
frame_a = tk.Frame()
label_a = tk.Label(master=frame_a, text="I'm in Frame A")
label_a.pack()

frame_b = tk.Frame()
label_b = tk.Label(master=frame_b, text="I'm in Frame B")
label_b.pack()

# Order of `frame_a` and `frame_b` is swapped
frame_b.pack()
frame_a.pack()

window.mainloop()
```

The output looks like this:



Now `label_b` is on top. Since `label_b` was assigned to `frame_b`, it moves to wherever `frame_b` is positioned.

All four of the widget types you have learned about — `Label`, `Button`, `Entry`, and `Text` — have a `master` attribute that is set when you instantiate them. That way you can control which `Frame` a widget is assigned to.

Frame widgets are great for organizing other widgets in a logical manner. Related widgets can be assigned to the same frame so that if the frame is ever moved in the window, the related widgets stay together.

In addition to grouping your widgets logically, Frame widgets can add a little flare to the visual presentation of your application. Read on to see how to create various borders for Frame widgets.

Adjusting Frame Appearance With Reliefs

Frame widgets can be configured with a `relief` attribute that creates a border around the frame. You can set `relief` to be any of the following values:

- `tk.FLAT`, no border effect (this is the default value).
- `tk.SUNKEN`, which creates a sunken effect.
- `tk.RAISED`, which creates a raised effect.
- `tk.GROOVE`, which creates a grooved border effect.
- `tk.RIDGE`, which creates a ridged effect.

To apply the border effect, you must set the `borderwidth` attribute to a value greater than 1. This attribute adjusts the width of the border, in pixels.

The best way to get a feel for what each effect looks like is to see them for yourself. Here is a script that packs five Frame widgets into a window, each with a different value for the `relief` argument.

```
import tkinter as tk

# 1
border_effects = {
    "flat": tk.FLAT,
    "sunken": tk.SUNKEN,
    "raised": tk.RAISED,
    "groove": tk.GROOVE,
```

```
    "ridge": tk.RIDGE,
}

window = tk.Tk()

for relief_name, relief in border_effects.items():
    # 2
    frame = tk.Frame(master=window, relief=relief, borderwidth=5)
    # 3
    frame.pack(side=tk.LEFT)
    # 4
    label = tk.Label(master=frame, text=relief_name)
    label.pack()

window.mainloop()
```

Let's break that script down some.

First, a dictionary is created whose keys are the names of the different relief effects available in Tkinter, and whose values are the corresponding Tkinter objects. This dictionary is assigned to the `border_effects` variable (#1).

After creating the `window` object, a `for` loop is used to loop over each item in the `border_effects` dictionary. At each step in the loop:

- A new `Frame` widget is created and assigned to the `window` object (#2). The `relief` attribute is set to the corresponding relief in the `border_effects` dictionary, and the `border` attribute is set to 5 so that the effect is visible.
- The `Frame` is then packed into the `window` using the `.pack()` method (#3). The `side` keyword argument you see is telling Tkinter which direction to pack the `frame` objects. You'll see more on how this works in the next section.
- A `Label` widget is created to display the name of the relief and is packed into the `frame` object just created (#4).

The window produced by the above script looks like this:



In this image, you can see that:

- `tk.FLAT` creates a flat looking frame
- `tk.SUNKEN` adds a border that gives the frame the appearance of being sunk into the window
- `tk.RAISED` gives the frame a border that makes it appear to protrude from the screen
- `tk.GROOVE` adds a border that appears as a sunken groove around an otherwise flat frame
- `tk.RIDGE` gives the appearance of a raised lip around the edge of the frame

Widget Naming Conventions

When you create a widget you can give it any name you like as long as it is a valid Python identifier. It is usually a good idea, though, to include the name of the widget class in the variable name you assign to the widget instance.

For example, if a `Label` widget is used to display a user's name, you might name the widget `label_user_name`. An `Entry` widget used to collect a user's age might be called `entry_age`.

When you include the widget class name in the variable name, you help yourself and anyone else that needs to read your code understand what type of widget to which the variable name refers.

Using the full name of the widget class can lead to long variable names, so you may want to adopt a shorthand for referring to each widget type. For the rest of this chapter, we'll use the following shorthand prefixes to name widgets:

Widget Class	Variable Name Prefix	Example
Label	lbl	lbl_name
Button	btn	btn_submit
Entry	ent	ent_age
Text	txt	txt_notes
Frame	frm	frm_address

In this section, you learned how to create a window, use widgets, and work with frames. At this point, you can make some simple windows displaying some messages, but a full-blown application is still out of reach.

In the next section, you'll learn how to control the layout of your applications using Tkinter's powerful geometry managers.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Try to re-create all of the screenshots in this section without looking at the source code. If you get stuck, check the code and finish your re-creation. Then wait for 10–15 minutes and try again.

Repeat this until you can produce all of the screenshots on your own. Focus on the output. It's okay if your own code is slightly different from the code in the book.

2. Write a program that displays a `Button` widget that is 50 text units wide and 25 text units tall and has a white background with blue text that reads "Click here".

3. Write a program that displays an `Entry` widget that is 40 text units wide and has a white background and black text. Use the `.insert()` method to display text in it that reads "What is your name?".

[Leave feedback on this section »](#)

18.6 Controlling Layout With Geometry Managers

Up until now, you've been adding widgets to windows and `Frame` widgets using the `.pack()` method, but you haven't been told what exactly this method does. Let's clear things up!

Application layout in Tkinter is controlled with **geometry managers**. `.pack()` is an example of a geometry manager, but it isn't the only one. Tkinter has two others: `.place()` and `.grid()`.

Each window and `Frame` in your application can use only one geometry manager. However, different frames can use different geometry managers, even if they are assigned to a frame or window using another geometry manager.

Let's start by taking a closer look at `.pack()`.

The `.pack()` Geometry Manager

`.pack()` uses a **packing algorithm** to place widgets in a `Frame` or window in a specified order. The packing algorithm has two primary steps. For a given widget, the algorithm:

1. Computes a rectangular area, called a **parcel**, that is just tall (or wide) enough to hold the widget and fills the remaining width (or height) in the window with blank space.
2. Centers the widget in the parcel, unless a different location is specified.

.pack() is powerful, but can be difficult to visualize. The best way to get a feel for .pack() is to look at some examples.

Let's see what happens when you .pack() three Label widgets into a Frame:

```
import tkinter as tk

window = tk.Tk()

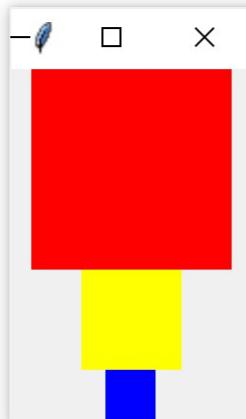
frame1 = tk.Frame(master=window, width=100, height=100, bg="red")
frame1.pack()

frame2 = tk.Frame(master=window, width=50, height=50, bg="yellow")
frame2.pack()

frame3 = tk.Frame(master=window, width=25, height=25, bg="blue")
frame3.pack()

window.mainloop()
```

By default, .pack() places each Frame below the previous one, in the order that they are assigned to the window:



Each `Frame` is placed at the top-most available position. The red `Frame` is placed at the top of the window. Then the yellow `Frame` is placed just below the red one, and the blue `Frame` just below the yellow one.

There are three invisible parcels containing each of the three `Frame` widgets. Each parcel is as wide as the window and as tall as the `Frame` that it contains. Since no anchor point was specified when `.pack()` was called for each `Frame`, they are all centered inside of their parcels. That's why each `Frame` is centered in the window.

`.pack()` accepts some keyword arguments for more precisely configuring widget placement. For example, you can set the `fill` keyword argument to specify which direction the frames should fill. The options are `tk.X` to fill in the horizontal direction, `tk.Y` to fill vertically, and `tk.BOTH` to fill in both directions.

Here's how you would stack the three frames so that each one fills the whole window horizontally:

```
import tkinter as tk

window = tk.Tk()

frame1 = tk.Frame(master=window, height=100, bg="red")
frame1.pack(fill=tk.X)

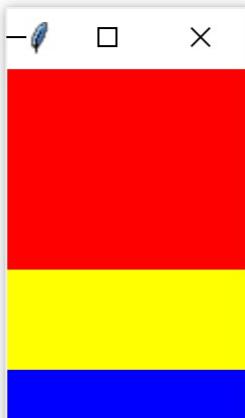
frame2 = tk.Frame(master=window, height=50, bg="yellow")
frame2.pack(fill=tk.X)

frame3 = tk.Frame(master=window, height=25, bg="blue")
frame3.pack(fill=tk.X)

window.mainloop()
```

Notice that the `width` is not set on any of the `Frame` widgets. `width` is no longer necessary because the `.pack()` method on each frame is set to fill horizontally, overriding any `width` you may set.

The window produced by this script looks like this:



One of the nice things about filling the window with `.pack()` is that the fill is responsive to window resizing. Try widening the window generated by the previous script to see how this works.

As you widen the window, the width of the three `Frame` widgets grow to fill the window. Notice, though, that the `Frame` widgets do not expand in the vertical direction.

The `side` keyword argument of `.pack()` specifies on which side of the window the widget should be placed. The available options are `tk.TOP`, `tk.BOTTOM`, `tk.LEFT`, and `tk.RIGHT`. If you do not set `side`, `.pack()` automatically used `tk.TOP` and places new widgets at the top of the window, or at the top-most portion of the window that isn't already occupied by a widget.

For example, the following script places three frames side by side from left to right and expands each frame to fill the window vertically:

```
import tkinter as tk
```

```
window = tk.Tk()

frame1 = tk.Frame(master=window, width=200, height=100, bg="red")
frame1.pack(fill=tk.Y, side=tk.LEFT)

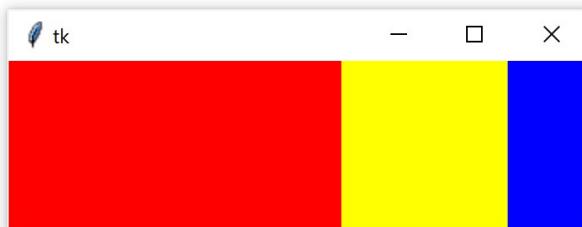
frame2 = tk.Frame(master=window, width=100, bg="yellow")
frame2.pack(fill=tk.Y, side=tk.LEFT)

frame3 = tk.Frame(master=window, width=50, bg="blue")
frame3.pack(fill=tk.Y, side=tk.LEFT)

window.mainloop()
```

This time, you have to specify the `height` keyword argument on at least one of the frames to force the window to have some height.

The resulting window looks like this:



Just like setting `fill=tk.X` made the frames resize responsively when the window is resized horizontally, setting `fill=tk.Y` makes the frames resize responsively when the window is resized vertically. Try it out!

To make the layout truly responsive, you can set an initial size for your frames using the `width` and `height` attributes. Then set the `fill` keyword argument of the `.pack()` method to `tk.BOTH` and set the `expand` keyword argument to `True`:

```
import tkinter as tk

window = tk.Tk()

frame1 = tk.Frame(master=window, width=200, height=100, bg="red")
frame1.pack(fill=tk.BOTH, side=tk.LEFT, expand=True)

frame2 = tk.Frame(master=window, width=100, bg="yellow")
frame2.pack(fill=tk.BOTH, side=tk.LEFT, expand=True)

frame3 = tk.Frame(master=window, width=50, bg="blue")
frame3.pack(fill=tk.BOTH, side=tk.LEFT, expand=True)

window.mainloop()
```

When you run the above script, you see a window that initially looks the same as the one generated in the previous example. The difference is that now you can resize the window however you want and the frames expand and fill the window responsively. Pretty cool!

The `.place()` Geometry Manager

You can use the `.place()` method of a widget to control the precise location that it should occupy in a window or `Frame`. You must provide two keyword arguments `x` and `y` that specify the x- and y-coordinates for the top-left corner of the widget. Both `x` and `y` are measured in pixels, not text units.

Keep in mind that the origin, where `x` and `y` are both 0, is the top left corner of the `Frame` or window, so you can think of the `y` argument of `.place()` as the number of pixels from the top of the window, and the `x` argument as the number of pixels from the left of the window.

Here's an example of how the `.place()` geometry manager works:

```
import tkinter as tk
```

```
window = tk.Tk()

# 1
frame = tk.Frame(master=window, width=150, height=150)
frame.pack()

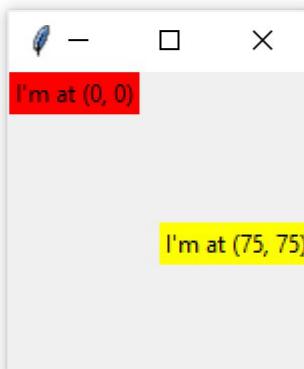
# 2
label1 = tk.Label(master=frame, text="I'm at (0, 0)", bg="red")
label1.place(x=0, y=0)

# 3
label2 = tk.Label(master=frame, text="I'm at (75, 75)", bg="yellow")
label2.place(x=75, y=75)

window.mainloop()
```

First, you create a new `Frame` widget called `frame1` that is 150 pixels wide and 150 pixels tall, and pack it into the window with `.pack()` (#1). Then you create a new `Label` called `label1` with a yellow background (#2) and place it in `frame1` at position (0, 0). Finally, you create a second `Label` called `label2` with a red background (#3) and place it in `frame1` at position (75, 75).

Here's the window the code produces:



.place() is not used often. It has two main drawbacks:

1. **Layout can be difficult to manage with .place()**, especially if your application has lots of widgets.
2. **Layouts created with .place() are not responsive**. They do not change as the window is resized.

One of the main challenges of cross-platform GUI development is making layouts that look good no matter which platform they are viewed on. .place() is a poor choice for making responsive and cross-platform layouts.

That's not to say .place() should never be used. It might be just what you need. For example, if you are creating a GUI interface for a map, then .place() might be the perfect choice to ensure widgets are placed at the correct distance from each other on the map.

.pack() is usually a better choice than .place(), but even .pack() has some downsides. The placement of widgets depends on the order in which .pack() is called, so it can be difficult to modify existing applications without fully understanding the code controlling the layout.

The .grid() geometry manager solves a lot of these issues, as you'll see in the next section.

The .grid() Geometry Manager

The geometry manager you will likely use most often is the .grid() geometry manager. .grid() provides all the power of .pack() in a format that is easier to understand and maintain.

.grid() works by splitting a window or `Frame` into rows and columns. You specify the location of a widget by calling .grid() and passing the row and column indices to the `row` and `column` keyword argument, respectively. Both row and column indices start at 0, so a row index of 1 and a column index of 2 tells .grid() to place a widget in the third column of the second row.

For example, the following script creates a 3×3 grid of frames with Label widgets packed into them:

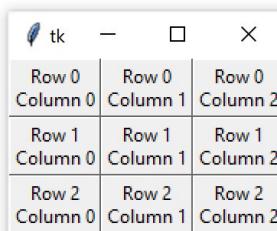
```
import tkinter as tk

window = tk.Tk()

for i in range(3):
    for j in range(3):
        frame = tk.Frame(
            master=window,
            relief=tk.RAISED,
            borderwidth=1
        )
        frame.grid(row=i, column=j)
        label = tk.Label(master=frame, text=f"Row {i}\nColumn {j}")
        label.pack()

window.mainloop()
```

Here's what the resulting window looks like:



Two geometries managers are being used in this example. Each `Frame` is attached to the `window` with the `.grid()` geometry manager, and each `label` is attached to its master `Frame` with `.pack()`.

The important thing to realize here is that even though `.grid()` is called on each `Frame` object, the geometry manager applies to the `window` object. Similarly, the layout of each `frame` is controlled with the `.pack()` geometry manager.

The frames in the previous example are placed tightly next to one another. To add some space around each `Frame`, you can set the padding of each cell in the grid. **Padding** is just some blank space that surrounds a widget and separates it visually from its contents.

There are two types of padding: **external padding** and **internal padding**. External padding adds some space around the outside of a grid cell. It is controlled with two keyword arguments of `.grid()`:

1. `padx`, which adds padding in the horizontal direction
2. `pady`, which adds padding in the vertical direction.

Both `padx` and `pady` are measured in pixels, not text units, so setting both of them to the same value will create the same amount of padding in both directions.

Let's add some padding around the outside of the frames in the previous example:

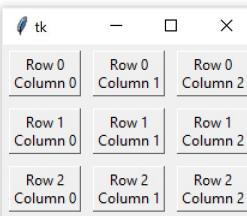
```
import tkinter as tk

window = tk.Tk()

for i in range(3):
    for j in range(3):
        frame = tk.Frame(
            master=window,
            relief=tk.RAISED,
            borderwidth=1
        )
        frame.grid(row=i, column=j, padx=5, pady=5)
        label = tk.Label(master=frame, text=f"Row {i}\nColumn {j}")
```

```
label.pack()  
  
window.mainloop()
```

Here's the resulting window:

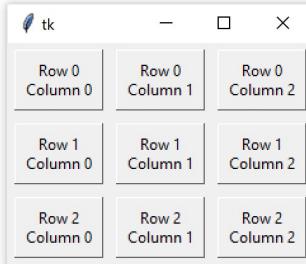


.pack() also has `padx` and `pady` parameters. The following code is nearly identical to the previous code, except that 5 pixels of additional padding have been added around each Label in the both the x and y directions:

```
import tkinter as tk  
  
window = tk.Tk()  
  
for i in range(3):  
    for j in range(3):  
        frame = tk.Frame(  
            master=window,  
            relief=tk.RAISED,  
            borderwidth=1  
        )  
        frame.grid(row=i, column=j, padx=5, pady=5)  
  
        label = tk.Label(master=frame, text=f"Row {i}\nColumn {j}")  
        label.pack(padx=5, pady=5)
```

```
window.mainloop()
```

The extra padding around the `Label` widgets gives each cell in the grid a little bit of breathing room between the `Frame` border and the text in the `Label`:



That looks pretty nice! But, if you try and expand the window in any direction, you'll notice that the layout isn't very responsive. The whole grid stays at the top left corner as the window expands.

You can adjust how the rows and columns of the grid grow as the window is resized using the `.columnconfigure()` and `.rowconfigure()` methods on the `window` object. Remember, the grid is attached to `window`, even though you are calling `.grid()` on each `Frame` widget.

Both `.columnconfigure()` and `.rowconfigure()` take three essential arguments:

1. The index of the grid column or row that you want to configure. You may also specify a list of indices to configure multiple rows or columns at the same time.
2. A keyword argument called `weight` that determines how the column or row should respond to window resizing relative to the other columns and rows.
3. A keyword argument called `minsize` that sets the minimum size of the row height or column width in pixels.

`weight` is set to 0 by default, which means that the column or row does not expand as the window resizes. If every column and row is given a weight of 1, they all grow at the same rate. If one column has a weight of 1 and another a weight of 2, then the second column expands at twice the rate of the first.

Let's adjust the previous script to better handle window resizing:

```
import tkinter as tk

window = tk.Tk()

for i in range(3):
    window.columnconfigure(i, weight=1, minsize=75)
    window.rowconfigure(i, weight=1, minsize=50)

    for j in range(0, 3):
        frame = tk.Frame(
            master=window,
            relief=tk.RAISED,
            borderwidth=1
        )
        frame.grid(row=i, column=j, padx=5, pady=5)

        label = tk.Label(master=frame, text=f"Row {i}\nColumn {j}")
        label.pack(padx=5, pady=5)

window.mainloop()
```

The `.columnconfigure()` and `.rowconfigure()` methods are placed in the body of the outer `for` loop. You could explicitly configure each column and row outside of the `for` loop, but that would require writing an additional six lines of code.

On each iteration of the loop, the `i`-th column and row are configured to have a `weight` of 1. This ensures that each row and column expands at the same rate whenever the window is resized.

The `minsize` argument is set to 75 for each column and 50 for each row. This makes sure the `Label` widget always displays its text without chopping off any characters, even if the window size is extremely small.

Try running the script to get a feel for how it works! Play around with the `weight` and `minsize` parameters to see how they affect the grid.

By default, widgets are centered in their grid cells. For example, the following code creates two `Label` widgets and places them in a grid with one column and two rows:

```
import tkinter as tk

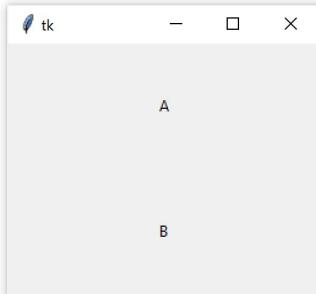
window = tk.Tk()
window.columnconfigure(0, minsize=250)
window.rowconfigure([0, 1], minsize=100)

label1 = tk.Label(text="A")
label1.grid(row=0, column=0)

label2 = tk.Label(text="B")
label2.grid(row=1, column=0)

window.mainloop()
```

Each grid cell is 250 pixels wide and 100 pixels tall. The labels are placed in the center of each cell, as you can see in the following figure:



You can change the location of each label inside of the grid cell using the `.grid()` method's `sticky` parameter. `sticky` accepts a string containing one or more of the following letters:

- "n" or "N" to align to the top center of the cell
- "e" or "E" to align to the right center side of the cell
- "s" or "S" to align to the bottom center of the cell
- "w" or "W" to align to the left center side of the cell

The letters "n", "e", "s", and "w" come from the cardinal directions, north, south, east, and west.

For example, setting `sticky` to "n" on both `Labels` in the previous code positions each `Label` at the top center of its grid cell:

```
import tkinter as tk

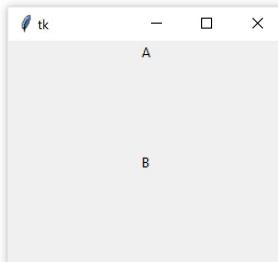
window = tk.Tk()
window.columnconfigure(0, minsize=250)
window.rowconfigure([0, 1], minsize=100)

label1 = tk.Label(text="A")
label1.grid(row=0, column=0, sticky="n")

label2 = tk.Label(text="B")
```

```
label2.grid(row=1, column=0, sticky="n")  
  
window.mainloop()
```

Here's the output:

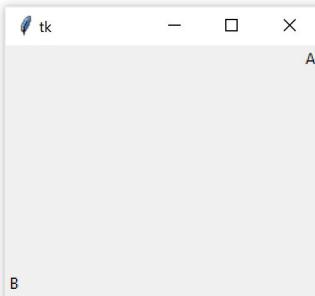


You can combine multiple letters in a single string to position each Label in a corner of its grid cell:

```
import tkinter as tk  
  
window = tk.Tk()  
window.columnconfigure(0, minsize=250)  
window.rowconfigure([0, 1], minsize=100)  
  
label1 = tk.Label(text="A")  
label1.grid(row=0, column=0, sticky="ne")  
  
label2 = tk.Label(text="B")  
label2.grid(row=1, column=0, sticky="sw")  
  
window.mainloop()
```

In this example, the `sticky` parameter of `label1` is set to "ne", which places the label at the top right corner of its grid cell. `label2` is positioned in the bottom left corner by passing "sw" to `sticky`.

Here's what that looks like in the window:



When a widget is positioned with `sticky`, the size of the widget itself is just big enough to contain any text and other contents inside of it. It won't fill the entire grid cell.

In order to fill the grid, you can specify "ns", which forces the widget to fill the cell in the vertical direction, or "ew" to fill the cell in the horizontal direction. To fill the entire cell, set `sticky` to "nsew".

The following example illustrates each of these options:

```
import tkinter as tk

window = tk.Tk()

window.rowconfigure(0, minsize=50)
window.columnconfigure([0, 1, 2, 3], minsize=50)

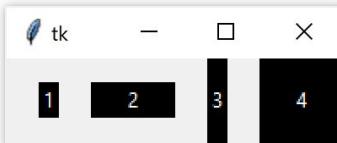
label1 = tk.Label(text="1", bg="black", fg="white")
label2 = tk.Label(text="2", bg="black", fg="white")
label3 = tk.Label(text="3", bg="black", fg="white")
label4 = tk.Label(text="4", bg="black", fg="white")

label1.grid(row=0, column=0)
label2.grid(row=0, column=1, sticky="ew")
```

```
label3.grid(row=0, column=2, sticky="ns")
label4.grid(row=0, column=3, sticky="nsew")

window.mainloop()
```

Here's what the output looks like:



What the above example illustrates is that the `.grid()` geometry manager's `sticky` parameter can be used to achieve the same effects as the `.pack()` geometry manager's `fill` parameter.

The correspondence between the `sticky` and `fill` parameters is summarized in the following table:

<code>.grid()</code>	<code>.pack()</code>
<code>sticky="ns"</code>	<code>fill=tk.Y</code>
<code>sticky="ew"</code>	<code>fill=tk.X</code>
<code>sticky="nsew"</code>	<code>fill=tk.BOTH</code>

`.grid()` is a powerful geometry manager. It is often easier to understand than `.pack()` and is much more flexible than `.place()`. When creating new Tkinter applications, consider using `.grid()` as your primary geometry manager.

Note

.grid() offers much more flexibility than you have seen here. For example, you can configure cells to span multiple rows and columns.

For more information, check out the [Grid Geometry Manager](#) section of the [TkDocs tutorial](#).

Now that you've got the basics of Tkinter's geometry managers down, the next step is to assign actions to buttons to bring your applications to life.

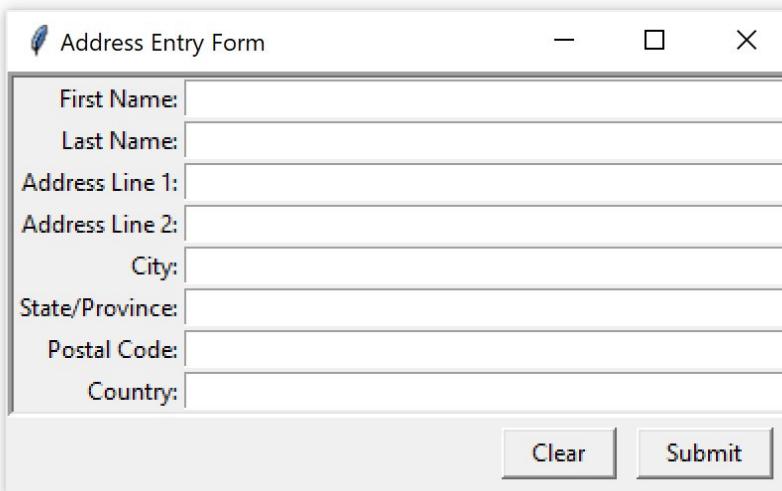
Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Re-create all of the screenshots in this section without looking at the source code. If you get stuck, check the code and finish your re-creation. Then wait for 10–15 minutes and try again. Repeat this until you can produce all of the screenshots on your own.

Focus on the output. It's okay if your own code is slightly different from the code in the book.)

2. Below is an image of a window made with Tkinter. Try and re-create the window using the techniques you have learned thus far. You may use any geometry manager you like.



[Leave feedback on this section »](#)

18.7 Making Your Applications Interactive

By now, you have a pretty good idea how to create a window with Tkinter, add some widgets, and control the application layout. That's great! But applications shouldn't just look good. They need to actually do something!

In this section, you'll learn how to bring your applications to life by performing actions whenever certain events occur.

Events and Event Handlers

When you create a Tkinter application, you must call the `window.mainloop()` method to start the **event loop**. During the event loop, your application checks if an event has occurred. If so, then some code can be executed in response.

The event loop is provided for you with Tkinter, so you do not have to write any code that checks for event yourself. However, you do have to write the code that is executed in response to an event. In Tkinter, you write functions called **event handlers** for the events that you use in your application.

So what is an event, and what happens when one occurs? An **event** is any action that occurs during the event loop that might trigger some behavior in the application, such as when a key or mouse button is pressed.

When an event occurs, an **event object** is emitted, which means that an instance of a class representing the event is instantiated. You don't need to worry about creating these classes yourself. Tkinter will create instances of event classes for you automatically.

Let's write our own event loop in order to better understand how Tkinter's event loop works. That way you can see how Tkinter's event loop fits into your application, and which parts you need to write yourself.

Assume there's a list called `events_list` that contains event objects. A new event object is automatically appended to `events_list` every time an event occurs in your program. We don't need to implement this updating mechanism. It just magically happens for us in this make-believe example.

Using an infinite loop, we can continually check if there are any event objects in `events_list`:

```
# Assume that this list gets updated automatically
events_list = []

# Run the event loop
while True:
    # If events_list is empty, the no events have occurred and we
    # can skip to the next iteration of the loop
    if events_list == []:
```

```
continue
```

```
# If execution reaches this point, then there is at least one
# event object in events_list
event = events_list[0]
```

Right now, the event loop we have created doesn't do anything with event. Let's change that.

Suppose our application needs to respond to key presses. We need to check that event was generated by a user pressing a key on their keyboard, and, if so, pass event to an event handler function for key presses.

We'll assume that event has a `.type` attribute set to the string "keypress" if the event is a keypress event object, and a `.char` attribute containing the character of the key that was pressed.

Note

Let's add a `handle_keypress()` function and update our event loop code:

```
events_list = []

# Create an event handler
def handle_keypress(event):
    """Print the character associated to the key pressed"""
    print(event.char)

while True:
    if events_list == []:
        continue
    event = events_list[0]

    # If event is a keypress event object
    if event.type == "keypress":
```

```
# Call the keypress event handler
handle_keypress(event)
```

When you call Tkinter's `window.mainloop()` method, something like the above loop is run for you! `.mainloop()` takes care of two parts of the loop for you: it maintains a list of events that have occurred, and it runs an event handler any time a new event is added to that list.

Let's update our event loop to use `window.mainloop()` instead of our own event loop:

```
import tkinter as tk

# Create a window object
window = tk.Tk()

# Create an event handler
def handle_keypress(event):
    """Print the character associated to the key pressed"""
    print(event.char)

# Run the event loop
window.mainloop()
```

`.mainloop()` takes care of a lot for us, but there's something missing from the above code. How does Tkinter know when to use `handle_keypress()`? Tkinter widgets have a `.bind()` method to do just this.

The `.bind()` Method

To call an event handler whenever an event occurs on a widgets, use the widget's `.bind()` method. The event handler is said to be **bound** to the event, because it is called every time the event occurs.

Continuing with the keypress example from the previous section, let's use `.bind()` to bind `handle_keypress()` to the keypress event:

```
import tkinter as tk

window = tk.Tk()

def handle_keypress(event):
    """Print the character associated to the key pressed"""
    print(event.char)

# Bind keypress event to handle_keypress()
window.bind("<Key>", handle_keypress)

window.mainloop()
```

Here, the `handle_keypress()` event handler is bound to a "`<Key>`" event using the `window.bind()` method. Whenever a key is pressed while the application is running, the character of the key pressed will be printed.

`.bind()` always takes two arguments:

1. An event, which is represent with a string of the form "`<event_name>`", where `event_name` can be any of Tkinter's events.
2. An event handler, which is the name of the function to be called whenever the event occurs.

The event handler is bound to the widget on which `.bind()` is called. When the event handler is called, the event object is passed to the event handler function.

In the example above, the event handler is bound to the window itself, but you can bind an event handler to any widget in your application.

For example, you can bind an event handler to a `Button` widget that will perform some action whenever the button is pressed:

```
def handle_click(event):
    print("The button was clicked!")
```

```
button = tk.Button(text="Click me!")  
  
button.bind("<Button-1>", handle_click)
```

In this example, the "<Button-1>" event on the `button` widget is bound to the `handle_click` event handler. The "<Button-1>" event occurs whenever the left mouse button is pressed while the mouse is over the widget.

There are other events for mouse button clicks including "<Button-2>" for the middle mouse button, if one exists, and "<Button-3>" for the right mouse button.

Note

For a list of commonly used events, see the [Event types](#) section of the [Tkinter 8.5 reference](#).

You can bind any event handler to any kind of widget with `.bind()`, but there is an easier way to bind event handlers to button clicks using the `Button` widget's `command` attribute.

The `command` Attribute

Every `Button` widget has a `command` attribute that you can assign to a function. Whenever the button is pressed, the function is executed.

Let's look at an example. First, we'll create a window with a `Label` widget that holds a numerical value. We'll put two buttons on the left and right of the label. The left button will be used to decrease the value in the `Label`, and the right one will increase the value.

Here's the code for the window:

```
import tkinter as tk  
  
window = tk.Tk()
```

```
window.rowconfigure(0, minsize=50, weight=1)
window.columnconfigure([0, 1, 2], minsize=50, weight=1)

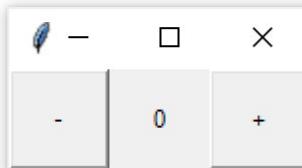
btn_decrease = tk.Button(master=window, text="-")
btn_decrease.grid(row=0, column=0, sticky="nsew")

lbl_value = tk.Label(master=window, text="0")
lbl_value.grid(row=0, column=1)

btn_increase = tk.Button(master=window, text="+")
btn_increase.grid(row=0, column=2, sticky="nsew")

window.mainloop()
```

The window looks like this:



With the app layout defined, let's bring it to life by giving the buttons some commands.

Let's start with the left button. When this button is pressed it should decrease the value in the label by 1. There are two things we need to know how to do in order to do this: how to get the text in a `Label`, and how to update the text in a `Label`.

`Label` widgets don't have a `.get()` method like `Entry` and `Text` widgets do. However, you can retrieve the text from the label by accessing the `text` attribute with dictionary-style subscript notation:

```
label = Tk.Label(text="Hello")

# Retrieve a Label's text
text = label["text"]

# Set new text for the label
label["text"] = "Good bye"
```

Now that we know how to get and set a label's text, let's write a function `increase()` that increases the value in the `lbl_value` by 1:

```
def increase():
    value = int(lbl_value["text"])
    lbl_value["text"] = f"{value + 1}"
```

`increase()` gets the text from `lbl_value` and converts it to an integer with `int()`. Then it increases this value by 1 and sets the label's `text` attribute to this new value.

We also need a `decrease()` function that decreases the value in `lbl_value` by 1:

```
def decrease():
    value = int(lbl_value["text"])
    lbl_value["text"] = f"{value - 1}"
```

Put the `increase()` and `decrease()` functions in your code just after the `import` statement.

To connect the buttons to the functions, assign the function to the button's `command` attribute. You can do this when you instantiate the button. For example, to assign the `increase()` function to the `btn_increase`, update the line that instantiates the button to the following:

```
btn_increase = tk.Button(master=window, text="+", command=increase)
```

Now assign the `decrease()` function to `btn_decrease`:

```
btn_decrease = tk.Button(master=window, text="-", command=decrease)
```

That's all you need to do to bind the buttons to the `increase()` and `decrease()` functions and make the program functional. Try saving your changes and running the application!

Here's the full application code for your reference:

```
import tkinter as tk

def increase():
    value = int(lbl_value["text"])
    lbl_value["text"] = f"{value + 1}"

def decrease():
    value = int(lbl_value["text"])
    lbl_value["text"] = f"{value - 1}"

window = tk.Tk()

window.rowconfigure(0, minsize=50, weight=1)
window.columnconfigure([0, 1, 2], minsize=50, weight=1)

btn_decrease = tk.Button(master=window, text="-", command=decrease)
btn_decrease.grid(row=0, column=0, sticky="nsew")

lbl_value = tk.Label(master=window, text="0")
lbl_value.grid(row=0, column=1)

btn_increase = tk.Button(master=window, text="+", command=increase)
btn_increase.grid(row=0, column=2, sticky="nsew")

window.mainloop()
```

This app is not particularly useful, but the skills you learned here apply to every app you'll make:

- Use widgets to create the components of the user interface.
- Use geometry managers to control the layout of the application.
- Write functions that interact with various components to capture and transform user input.

In the next two sections, you'll build apps that do something useful. First, you'll build a temperature converter that converts a temperature input as Fahrenheit to Celsius. After that, you'll build a text editor that can open, edit and save text files!

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

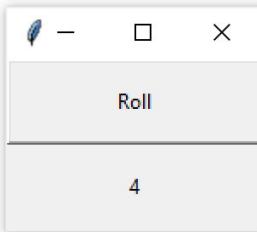
1. Write a program that displays a single button with the default background color and black text that reads "Click me".

When the user clicks on the button, the button background should change to a color randomly selected from the following list:

```
["red", "orange", "yellow", "blue", "green", "indigo", "violet"]
```

2. Write a program that simulates rolling a six-sided die. There should be one button with the text "Roll". When the user clicks the button, a random integer from 1 to 6 should be displayed.

The application window should look something like this:



[Leave feedback on this section »](#)

18.8 Example App: Temperature Converter

In this section, you'll build a temperature converter that allows the user to input a temperature in degrees Fahrenheit and push a button to convert that temperature to degrees Celsius.

We'll walk through the code step by step. You can also find the full source code at the end of this section for your reference.

To get the most out of this section, open up IDLE's script window and follow along.

Before we start coding, let's design the app. We need three basic elements:

1. An `Entry` widget called `ent_temperature` to enter the Fahrenheit value
2. A `Label` widget called `lbl_result` to display the Celsius result
3. A `Button` widget called `btn_convert` that reads the value from the `Entry` widget, converts it from Fahrenheit to Celsius, and sets the text of the `Label` widget to the result when clicked

We can arrange these in a grid with a single row and one column for each widget. That gets us a minimally working application, but it isn't very user friendly. Everything needs to have some helpful labels.

Let's put a label directly to the right of the `ent_temperature` widget containing the °F symbol so that the user knows that the value `ent_temperature` should be in degrees Fahrenheit. To do this, we'll set the label text to "\N{DEGREES FAHRENHEIT}", which uses Python's named Unicode character support to display the °F symbol.

We can give `btn_convert` a little flair by setting its text to the value "\N{RIGHTWARDS BLACK ARROW}", which displays a black arrow pointing to the right. We will also make sure that `lbl_result` always has the °C symbol "\N{DEGREES CELSIUS}" at the end to indicate that result is in degrees Celsius.

Here's what the final window will look like:



Now that we know what widgets we need and what they window is going to look like, let's start coding it up! First, import `tkinter` and create a new window:

```
import tkinter as tk

window = tk.Tk()
window.title("Temperature Converter")
```

`window.title()` sets the title of an existing window. When you finally run this application, the window will have the text *Temperature Converter* in its title bar.

Next, we'll create the `ent_temperature` widget with a label called `lbl_temp` and assign both of them to a `Frame` widget called `frm_entry`:

```
frm_entry = tk.Frame(master=window)
ent_temperature = tk.Entry(master=frm_entry, width=10)
lbl_temp = tk.Label(master=frm_entry, text="\N{DEGREE FAHRENHEIT}")
```

`ent_temperature` is where the user will enter the Fahrenheit value, and `lbl_temp` is used to label `ent_temperature` with the °F symbol. `frm_entry` is just a container that groups `ent_temperature` and `lbl_temp` together.

We want `lbl_temp` to be placed directly to the right of `ent_temperature`, so we can lay them out in the `frm_entry` using the `.grid()` geometry manager with one row and two columns:

```
ent_temperature.grid(row=0, column=0, sticky="e")
lbl_temp.grid(row=0, column=1, sticky="w")
```

We've set the `sticky` parameter to "e" for `ent_temperature` so that it always sticks to the right-most edge of its grid cell. Setting `sticky` to "w" for `lbl_temp` will keep it stuck to the left-most edge of its grid cell. This ensures that `lbl_temp` is always located immediately to the right of `ent_temperature`.

Now let's make the `btn_convert` and the `lbl_result` for converting the temperature entered into `ent_temperature` and displaying the results:

```
btn_convert = tk.Button(
    master=window,
    text="\N{RIGHTWARDS BLACK ARROW}"
)
lbl_result = tk.Label(master=window, text="\N{DEGREE CELSIUS}")
```

Like `frm_entry`, both `btn_convert` and `lbl_result` are assigned to `window`. Together, these three widgets make up the three cells in the main application grid. Let's use `.grid()` to go ahead and lay them out now:

```
frm_entry.grid(row=0, column=0, padx=10)
btn_convert.grid(row=0, column=1, pady=10)
lbl_result.grid(row=0, column=2, padx=10)
```

Finally, run the application:

```
window.mainloop()
```

That looks great, but the button doesn't do anything yet. At the top of your script file, just below the `import` line, add a function called `fahrenheit_to_celsius()`. This function reads the value from `ent_temperature`, converts it from Fahrenheit to Celsius, and then displays the result in `lbl_result`:

```
def fahrenheit_to_celsius():
    """Convert the value for Fahrenheit to Celsius and insert the
    result into lbl_result.
    """
    fahrenheit = ent_temperature.get()
    celsius = (5/9) * (float(fahrenheit) - 32)
    lbl_result["text"] = f'{round(celsius, 2)} \N{DEGREE CELSIUS}'
```

Now go down to the line where you define `btn_convert` and set its `command` parameter to `fahrenheit_to_celsius`:

```
btn_convert = tk.Button(
    master=window,
    text="\N{RIGHTWARDS BLACK ARROW}",
    command=fahrenheit_to_celsius  # <-- Add this line
)
```

That's it! You've created a fully functional temperature converter app in just 26 lines of code! Pretty cool, right?

Here's the full script for your reference:

```
import tkinter as tk

def fahrenheit_to_celsius():
    """Convert the value for Fahrenheit to Celsius and insert the
    result into lbl_result.
    """

```

```
fahrenheit = ent_temperature.get()
celsius = (5/9) * (float(fahrenheit) - 32)
lbl_result["text"] = f'{round(celsius, 2)} \N{DEGREE CELSIUS}"

# Set-up the window
window = tk.Tk()
window.title("Temperature Converter")
window.resizable(width=False, height=False)

# Create the Fahrenheit entry frame with an Entry
# widget and label in it
frm_entry = tk.Frame(master=window)
ent_temperature = tk.Entry(master=frm_entry, width=10)
lbl_temp = tk.Label(master=frm_entry, text="\N{DEGREE FAHRENHEIT}")

# Layout the temperature Entry and Label in frm_entry
# using the .grid() geometry manager
ent_temperature.grid(row=0, column=0, sticky="e")
lbl_temp.grid(row=0, column=1, sticky="w")

# Create the conversion Button and result display Label
btn_convert = tk.Button(
    master=window,
    text="\N{RIGHTWARDS BLACK ARROW}",
    command=fahrenheit_to_celsius
)
lbl_result = tk.Label(master=window, text="\N{DEGREE CELSIUS}")

# Set-up the layout using the .grid() geometry manager
frm_entry.grid(row=0, column=0, padx=10)
btn_convert.grid(row=0, column=1, pady=10)
lbl_result.grid(row=0, column=2, padx=10)

# Run the application
window.mainloop()
```

Let's take things up a notch. Read on to build a simple text editor.

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Try to re-create the temperature converter app without looking at the source code. If you get stuck, check the code and finish your re-creation. Then wait for 10–15 minutes and try again.

Repeat this until you can build the app from scratch on your own. Focus on the output. It's okay if your own code is slightly different from the code in the book.

[Leave feedback on this section »](#)

18.9 Example App: Text Editor

In this section, you'll build a text editor app that can create, open, edit, and save text files.

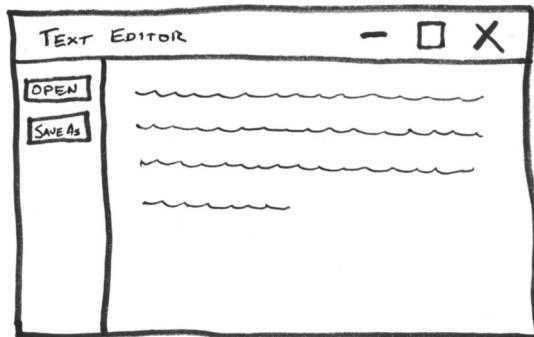
There are three essential elements in the application:

1. A `Button` widget called `btn_open` for opening a file for editing
2. A `Button` widget called `btn_save` for saving a file
3. A `TextBox` widget called `txt_edit` for creating and editing the text file

The three widgets will be arranged so that the two buttons are on the left-hand side of the window, and the text box is on the right-hand side.

The whole window should have a minimum height of 800 pixels, and `txt_edit` should have a minimum width of 800 pixels. The whole layout should be responsive so that if the window is resized then `txt_edit` is resized as well. The width of the `Frame` holding the buttons should not change, however.

Here's a sketch of how the window will look:



We can achieve the desired layout using the `.grid()` geometry manager. The layout contains a single row and two columns: a narrow column on the left for the buttons, and a wider column on the right for the text box.

To set the minimum sizes for the window and `txt_edit`, we can use the `minsize` parameters of the `.rowconfigure()` and `.columnconfigure()` window methods to 800. To handle resizing, we can set the `weight` parameters of these methods to 1.

In order to get both buttons into the same column, we'll need to create a `Frame` widget. Let's call that `Frame` widget `fr_buttons`. According to the sketch, the two buttons should be stacked vertically inside of this frame, with `btn_open` on top. We can do that with either the `.grid()` or `.pack()` geometry managers, but let's stick with `.grid()` since it is a little easier to work with.

Now that we have a plan, let's start coding the application. The first step is to create all of the widgets that we need:

```
import tkinter as tk

# 1
window = tk.Tk()
window.title("Simple Text Editor")
```

```
# 2
window.rowconfigure(0, minsize=800, weight=1)
window.columnconfigure(1, minsize=800, weight=1)

# 3
txt_edit = tk.Text(window)
fr_buttons = tk.Frame(window)
btn_open = tk.Button(fr_buttons, text="Open")
btn_save = tk.Button(fr_buttons, text="Save As...")
```

First we import `tkinter` and create a new window with the title "Simple Text Editor" (#1). Then the row and column configurations are set (#2). Finally, four widgets are created: the `txt_edit` text box, the `fr_buttons` frame, and the `btn_open` and `btn_save` buttons (#3).

Let's look at (#2) more closely. The `minsize` parameter of `.rowconfigure()` is set to 800 and `weight` is set to 1. The first argument is 0, so this sets the height of the first row to 800 pixels and makes sure that the height of the row grows proportionally to the height of the window. There is only one row in the application layout, so these settings apply to the entire window.

On the next line, `.columnconfigure()` is used to set the `width` and `weight` attributes of the column with index 1 to 800 and 1, respectively. Remember, row and column indices are zero based, so these settings apply only to the second column.

By configuring just the second column, the text box will expand and contract naturally when the window is resized while the column containing the buttons will remain at a fixed width.

Now we can work on the application layout. First, we'll assign the two buttons to the `fr_buttons` frame using the `.grid()` geometry manager:

```
btn_open.grid(row=0, column=0, sticky="ew", padx=5, pady=5)
btn_save.grid(row=1, column=0, sticky="ew", padx=5)
```

These two lines of code create a grid with two rows and one column in

the `fr_buttons` frame, since both `btn_open` and `btn_save` have their `master` attribute set to `fr_buttons`. `btn_open` is put in the first row and `btn_save` in the second row so that `btn_open` appears above `btn_save` in the layout, just we planned in our sketch.

Both `btn_open` and `btn_save` have their `sticky` attributes set to "ew", which forces the buttons to expand horizontally in both directions and fill the entire frame. This makes sure both buttons are the same size.

5 pixels of padding is placed around each button with the by setting the `padx` and `pady` parameters to 5. Only `btn_open` has vertical padding. Since it is on top, the vertical padding offsets the button down from the top of the window a bit and makes sure that there is a small gap between it and `btn_save`.

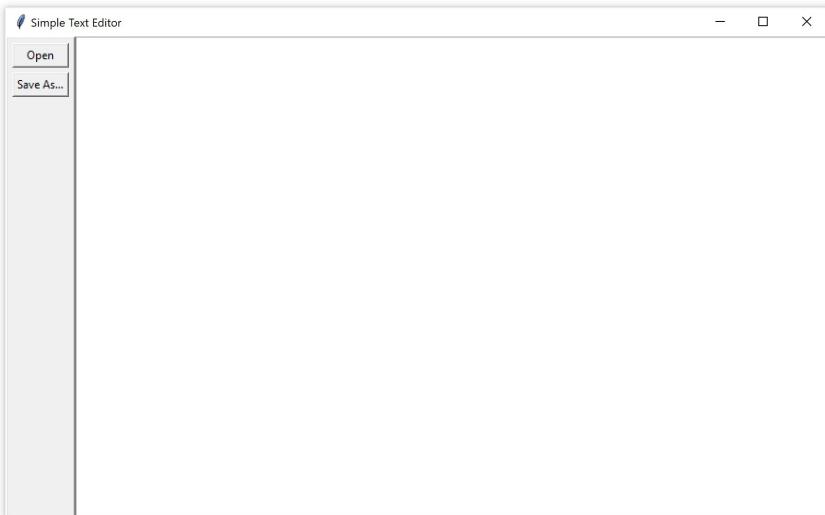
Now that `fr_buttons` is laid out and ready to go we can set up the grid layout for the rest of the window:

```
fr_buttons.grid(row=0, column=0, sticky="ns")
txt_edit.grid(row=0, column=1, sticky="nsew")
```

These two lines of code create a grid with one row and two columns for `window`. `fr_buttons` is placed in the first column and `txt_edit` in the second column so that `fr_buttons` appears to the left of `txt_edit` in the window layout.

The `sticky` parameter for `fr_buttons` is set to "ns", which forces the whole frame to expand vertically and fill the entire height of its column. `txt_edit` fills its entire grid cell because its `sticky` parameters is set to "nsew", which forces it to expand in every direction.

Now that the application lay out is complete, add `window.mainloop()` to the bottom of the program and save and run the file. The following window is displayed:



That looks great! But it doesn't do anything yet, so let's start writing the commands for the buttons.

The `btn_open` button needs to show a file open dialog and allow the user to select a file. Then it needs to open that file and set the text of `txt_edit` to the contents of the file.

Here's a function `open_file()` that does just this:

```
def open_file():
    """Open a file for editing."""
    # 1
    filepath = askopenfilename(
        filetypes=[("Text Files", "*.txt"), ("All Files", "*.*")]
    )

    # 2
    if not filepath:
        return

    # 3
```

```
txt_edit.delete("1.0", tk.END)

# 4
with open(filepath, "r") as input_file:
    text = input_file.read()
    txt_edit.insert(tk.END, text)

# 5
window.title(f"Simple Text Editor - {filepath}")
```

The `askopenfilename` dialog from the `tkinter.filedialog` module is used to display a file open dialog and store the selected file path to the `filepath` variable (#1). If the user closes the dialog box or clicks the `Cancel` button, then `filepath` will be `None` and the function will return without executing any of the code to read the file and set the text of `txt_edit` (#2).

If the user does choose a file, though, then the current contents of `txt_edit` are cleared using the `.delete()` method (#3). Then the select file is opened and the contents of the file are read using the `.read()` method and stored as a string in the `text` variable. The string `text` is assigned to `txt_edit` using `.insert()` (#4).

Finally, the title of the window is set so that it contains the path of the open file (#5).

Now you can update the program so that `btn_open` calls `open_file()` whenever it is clicked. There are three things you need to do to update the program:

1. Import the `askopenfilename()` function from `tkinter.filedialog` by adding the following import to the top of your program:

```
from tkinter.filedialog import askopenfilename
```

2. Add the definition of `open_file()` just below the import statements.
3. Set the `command` attribute of `btn_opn` to `open_file`:

```
btn_open = tk.Button(fr_buttons, text="Open", command=open_file)
```

Save the file and run it to check that everything is working. Try opening a text file!

Note

If you have trouble getting the updates to work, you can skip ahead to the end of this section to see the full code for the text editor application.

With `btn_open` working, let's work on the function for `btn_save`. It needs to open a save file dialog box so that the user can choose where they would like to save the file. We'll use the `asksaveasfilename` dialog in the `tkinter.filedialog` module for this. It also needs to extract the text currently in `txt_edit` and write this to a file and the selected location.

Here's a function that does just this:

```
def save_file():
    """Save the current file as a new file."""
    # 1
    filepath = asksaveasfilename(
        defaultextension=".txt",
        filetypes=[("Text Files", "*.txt"), ("All Files", "*.*")],
    )

    # 2
    if not filepath:
        return

    # 3
    with open(filepath, "w") as output_file:
        text = txt_edit.get("1.0", tk.END)
        output_file.write(text)

    # 4
```

```
window.title(f"Simple Text Editor - {filepath}")
```

The `asksaveasfilename` dialog box is used to get the desired save location from the user. The selected file path is stored in the `filepath` variable (#1). If the user closes the dialog box or clicks the `Cancel` button, then `filepath` will be `None` and the function returns without executing any of the code to save the text to a file (#2).

If the user does select a file path, then a new file is created. The text from `txt_edit` is extracted with the `.get()` method and assigned to the variable `text` and written to the output file (#3).

Finally, the title of the window is updated so that the new file path is displayed in the window title (#4).

Now you can update the program so that `btn_save` calls `save_file()` when it is clicked. There are three things you need to do in order to update the program:

1. Import the `asksaveasfilename()` function from `tkinter.filedialog` by updating the import at the top of your script, like so:

```
from tkinter.filedialog import askopenfilename, asksaveasfilename
```

2. Add the definition of `save_file()` just below the `open_file()` definition.
3. Set the `command` attribute of `btn_save` to `save_file`:

```
btn_save = tk.Button(fr_buttons, text="Save As...", command=save_file)
```

Save the file and run it. You've now got a minimal, yet fully functional, text editor!

Here's the full script for your reference:

```
import tkinter as tk
from tkinter.filedialog import askopenfilename, asksaveasfilename
```

```
def open_file():
    """Open a file for editing."""
    filepath = askopenfilename(
        filetypes=[("Text Files", "*.txt"), ("All Files", "*.*")]
    )
    if not filepath:
        return
    txt_edit.delete(1.0, tk.END)
    with open(filepath, "r") as input_file:
        text = input_file.read()
        txt_edit.insert(tk.END, text)
    window.title(f"Simple Text Editor - {filepath}")

def save_file():
    """Save the current file as a new file."""
    filepath = asksaveasfilename(
        defaultextension="txt",
        filetypes=[("Text Files", "*.txt"), ("All Files", "*.*")],
    )
    if not filepath:
        return
    with open(filepath, "w") as output_file:
        text = txt_edit.get(1.0, tk.END)
        output_file.write(text)
    window.title(f"Simple Text Editor - {filepath}")

window = tk.Tk()
window.title("Simple Text Editor")
window.rowconfigure(0, minsize=800, weight=1)
window.columnconfigure(1, minsize=800, weight=1)

txt_edit = tk.Text(window)
fr_buttons = tk.Frame(window, relief=tk.RAISED, bd=2)
btn_open = tk.Button(fr_buttons, text="Open", command=open_file)
btn_save = tk.Button(fr_buttons, text="Save As...", command=save_file)
```

```
btn_open.grid(row=0, column=0, sticky="ew", padx=5, pady=5)
btn_save.grid(row=1, column=0, sticky="ew", padx=5)

fr_buttons.grid(row=0, column=0, sticky="ns")
txt_edit.grid(row=0, column=1, sticky="nsew")

window.mainloop()
```

You've now built two GUI applications in Python. In doing so, you've applied many of the topics you've learned about throughout this book. That's no small achievement, so take some time to feel good about what you've done!

You're now ready to tackle some applications on your own!

Review Exercises

You can find the solutions to these exercises and many other bonus resources online at realpython.com/python-basics/resources.

1. Try to re-create the text editor app without looking at the source code. If you get stuck, check the code and finish your re-creation. Then wait for 10–15 minutes and try again.

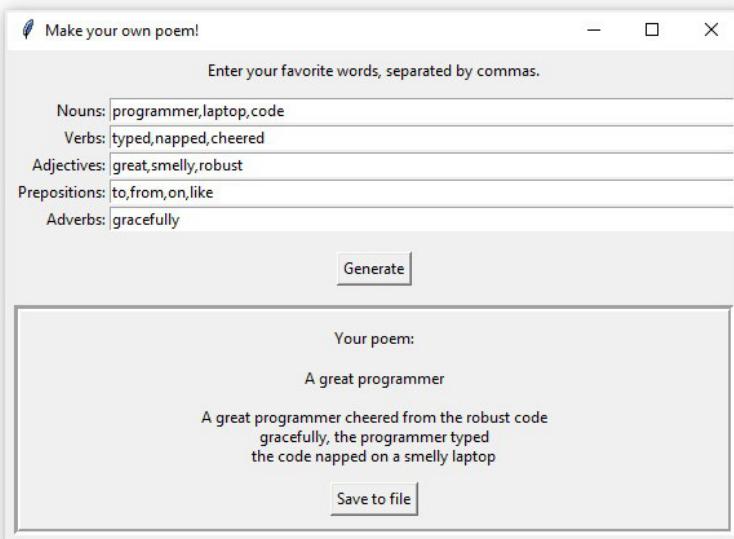
Repeat this until you can build the application from scratch on your own. Focus on the output. It's okay if your own code is slightly different from the code in the book.

[Leave feedback on this section »](#)

18.10 Challenge: Return of the Poet

For this challenge, you'll write a GUI application for generating poetry. This application is based off the poem generator from Chapter 9.

Visually, the application should look similar to this:



You may use whichever geometry manager you like, but the application should do all of the following:

1. Only allow the user to enter the correct number of words in each Entry widget:
 - At least 3 nouns
 - At least 3 verbs
 - At least 3 adjectives
 - At least 3 prepositions
 - At least 1 adverb

If too few words are entered into any of the Entry widgets, an error message should be displayed in the area where the generated poem is shown.

2. Randomly choose three nouns, adverbs, adjectives, and prepositions from the user input, and one adverb.
3. The program should generate the poem using the following template:

```
{A/An} {adj1} {noun1}  
  
A {adj1} {noun1} {verb1} {prep1} the {adj2} {noun2}  
{adverb1}, the {noun1} {verb2}  
the {noun2} {verb3} {prep2} a {adj3} {noun3}
```

4. The application must allow the user to export their poem to a file.
5. **Bonus:** Check that the user inputs unique words into each `Entry` widget. For example, if the user enters the same noun into the noun `Entry` widget twice, the application should display an error message when the user tries to generate the poem.

You can find the solutions to this code challenge and many other bonus resources online at realpython.com/python-basics/resources.

[Leave feedback on this section »](#)

18.11 Summary and Additional Resources

In this chapter, you learned how to build some simple graphical user interfaces (GUIs).

First, you learned how to use the EasyGUI package to create dialog boxes to display messages to a user, accept user input, and allow a user to select files for reading and writing. Then you learned about Tkinter, which is Python’s built-in GUI framework. Tkinter is more complex than EasyGUI, but also more flexible.

You learned how to work with widgets in Tkinter, including `Frame`, `Label`, `Button`, `Entry` and `Text` widgets. Widgets can be customized by assigning values to their various attributes. For example, setting the `text` attribute of a `Label` widget assigns some text to the label.

Next, you saw how to use Tkinter’s `.pack()`, `.place()` and `.grid()` geometry managers to give your GUI applications a layout. You learned how to control various aspects of the layout including internal and ex-

ternal padding, and how to create responsive layouts with the `.pack()` and `.grid()` managers.

Finally, you brought all of these skills together to create two full GUI applications: a temperature converter and a simple text editor.

Interactive Quiz

This chapter comes with a free online quiz to check your learning progress. You can access the quiz using your phone or computer at the following web address:

realpython.com/quizzes/python-basics-17

Additional Resources

To learn more about GUI programming in Python, check out these resources:

- [Tkinter tutorial](#)
- [Recommended resources on realpython.com](#)

[Leave feedback on this section »](#)

Chapter 19

Final Thoughts and Next Steps

Congratulations! You've made it to all the way to the end of this book. You already know enough to do a lot of amazing things with Python, but now the real fun starts: it's time to explore on your own!

The best way to learn is by solving real problems of your own. Sure, your code might not be very pretty or efficient when you're just starting out, but it will be useful. If you don't think you have any problems of the variety that Python could solve, pick a [popular module](#) that interests you and create your own project around it.

Part of what makes Python so great is the community. Know someone learning Python? Help them out! The only way to know you've really mastered a concept is when you can explain it to someone else.

Next up, dive into the more advanced material available at realpython.com or peruse the articles & tutorials featured in the PyCoder's Weekly newsletter.

When you feel ready, consider helping out with an open-source project on [GitHub](https://github.com). If puzzles are more your style, try working through some of the mathematical challenges on [Project Euler](https://projecteuler.net).

If you get stuck somewhere along the way, it's almost guaranteed that someone else has encountered (and potentially solved) the exact same problem before; search around for answers, particularly at [Stack Overflow](#), or find a [community](#) of Pythonistas willing to help you out.

If all else fails, `import this` and take a moment to meditate on that which is Python.

P.S. Come visit us on the web and continue your Python journey on the [realpython.com](#) website and the [@realpython](#) Twitter account.

19.1 Free Weekly Tips for Python Developers

Are you looking for a weekly dose of Python development tips to improve your productivity and streamline your workflows? Good news—we're running a free email newsletter for Python developers just like you.

The newsletter emails we send out are not your typical “here's a list of popular articles” flavor. Instead we aim for sharing at least one original thought per week in a (short) essay-style format.

If you'd like to see what all the fuss is about, then head on over to [realpython.com/newsletter](#) and enter your email address in the signup form. We're looking forward to meeting you!

19.2 Python Tricks: The Book

Now that you're familiar with the basics of Python, it's time to dig in deeper and to round out your knowledge.

With Real Python's *Python Tricks* book you'll discover Python's best practices and the power of beautiful & Pythonic code with simple examples and a step-by-step narrative.

You'll get one step closer to mastering Python, so you can write beautiful and idiomatic code that comes to you naturally.

Learning the ins and outs of Python is difficult—and with this book you'll be able to focus on taking your core Python skills to the next level.

Discover the “hidden gold” in Python’s standard library and start writing clean and Pythonic code today. Download a free sample chapter at realpython.com/pytricks-book

19.3 Real Python Video Course Library

Become a well-rounded Pythonista with Real Python’s large (and growing) collection of Python tutorials and in-depth training materials. With new content published weekly, you’ll always find something to boost your skills:

Master Practical, Real-World Python Skills: Our tutorials are created, curated, and vetted by a community of expert Pythonistas. At Real Python you’ll get the trusted resources you need on your path to Python mastery.

Meet Other Pythonistas: Join the Real Python Slack chat and meet the Real Python Team and other subscribers. Discuss your coding and career questions, vote on upcoming tutorial topics, or just hang out with us at this virtual water cooler.

Interactive Quizzes & Learning Paths: See where you stand and practice what you learn with interactive quizzes, hands-on coding challenges, and skills-focused learning paths.

Track Your Learning Progress: Mark lessons as completed or in-progress and learn at your own comfortable pace. Bookmark interesting lessons and review them later to boost long-term retention.

Completion Certificates: For each course you complete you receive a shareable (and printable) Certificate of Completion, hosted

privately on the Real Python website. Embed your certificates in your portfolio, LinkedIn resume, and other websites to show the world that you're a dedicated Pythonista.

Regularly Updated: Keep your skills fresh and keep up with technology. We're constantly releasing new members-only tutorials and update our content regularly.

See what's available at realpython.com/courses

19.4 PythonistaCafe: A Community for Python Developers

Mastering Python is *not* just about getting the books and courses to study. To be successful you also need a way to stay motivated and to grow your abilities in the long run.

Many Pythonistas we know are struggling with this. It's simply a lot less fun to build your Python skills completely alone.

If you're a self-taught developer with a non-technical day job, it's hard to grow your skills all by yourself. And with no coders in your personal peer group, there's nobody to encourage or support you in your endeavor of becoming a better developer.

Maybe you're already working as a developer, but no one else at your company shares your love for Python. It's frustrating when you can't share your learning progress with anyone or ask for advice when you feel stuck.

From personal experience, we know that existing online communities and social media don't do a great job at providing that support network either. Here are a few of the best, but they still leave a lot to be desired:

- *Stack Overflow* is for asking focused, one-off questions. It's hard to make a human connection with fellow commenters on the platform.

form. Everything is about the facts, not the people. For example, moderators will freely edit other people’s questions, answers, and comments. It feels more like a wiki than a forum.

- *Twitter* is like a virtual water cooler and great for “hanging out” but it’s limited to messages that can only be a few sentences long at a time—not great for discussing anything substantial. Also, if you’re not constantly online, you’ll miss out on most of the conversations. And if you *are* constantly online, your productivity takes a hit from the never-ending stream of interruptions and notifications. Slack chat groups suffer from the same flaws.
- *Hacker News* is for discussing and commenting on tech news. It doesn’t foster long-term relationships between commenters. It’s also one of the most aggressive communities in tech right now with little moderation and a borderline toxic culture.
- *Reddit* takes a broader stance and encourages more “human” discussions than Stack Overflow’s one-off Q&A format. But it’s a huge public forum with millions of users and has all of the associated problems: toxic behavior, overbearing negativity, people lashing out at each other, jealousy, ... In short, all the “best” parts of the human behavior spectrum.

Eventually I realized that what holds so many developers back is their limited access to the global Python coding community. That’s why I founded [PythonistaCafe](#), a peer-to-peer learning community for Python developers.



A good way to think of PythonistaCafe is to see it as a club of mutual improvement for Python enthusiasts:

Inside PythonistaCafe you'll interact with professional developers and hobbyists from all over the world who will share their experiences in a safe setting—so you can learn from them and avoid the same mistakes they've made.

Ask anything you want and it will remain private. You must have an active membership to read and write comments and as a paid community, trolling and offensive behavior are virtually nonexistent.

The people you meet on the inside are actively committed to improving their Python skills because membership in PythonistaCafe is invite-only. All prospective members are required to submit an application to make sure they're a good fit for the community.

You'll be involved in a community that understands you, and the skills and career you're building, and what you're trying to achieve. If you're trying to grow your Python skills but haven't found the support system you need, we're right there for you.

PythonistaCafe is built on a private forum platform where you can ask questions, get answers, and share your progress. We have members located all over the world and with a wide range of proficiency levels.

You can learn more about PythonistaCafe, our community values, and what we're all about at www.pythonistacafe.com.

19.5 Acknowledgements

This book would not have been possible without the help and support of so many friends and colleagues. We would like to thank many people for their assistance in making this book possible:

- **Our Families:** For bearing with us through “crunch mode” as we worked night and day to get this book into your hands.
- **The CPython Team:** For producing the amazing programming language and tools that we love and work with every day.
- **The Python Community:** For all the people who are working

hard to make Python the most beginner-friendly and welcoming programming language in the world, running conferences, and maintaining critical infrastructure like PyPI.

- **The Readers of realpython.com, Like You:** Thanks so much for reading our online articles and purchasing this book. Your continued support and readership is what makes all of this possible!

We hope that you will continue to be active in the community, asking questions and sharing tips. Reader feedback has shaped this book over the years and will continue to help us make improvements in future editions, so we look forward to hearing from you.

Finally, our deepest thanks to all of the Kickstarter backers who took a chance on this project in 2012. We never expected to gather such a large group of helpful, encouraging people.

[Leave feedback on this section »](#)

This is an Early Access version of “Python Basics: A Practical Introduction to Python 3”

With your help we can make this book even better:

At the end of each section of the book you’ll find a “magical” feedback link. Clicking the link takes you to an **online feedback form** where you can share your thoughts with us.

We welcome any and all feedback or suggestions for improvement you may have.

Please feel free to be as terse or detailed as you see fit. All feedback is stored anonymously, but you can choose to leave your name and contact information so we can follow up or mention you on our “Thank You” page.

We use a different feedback link for each section, so we’ll always know which part of the book your notes refer to.

Thank you for helping us make this book an even more valuable learning resource for the Python community.

— Dan Bader, Editor-in-Chief at Real Python
