# FOUNDATIONAL PYTHON FOR DATA SCIENCE



## KENNEDY BEHRMAN

# Foundational Python for Data Science

*This page intentionally left blank*

# Foundational Python for Data Science

Kennedy R. Behrman

✦✦ Addison-Wesley

❖

*This book is dedicated to Tatiana, Itta, and Maple,*
*who is probably still under the bed.*

❖

# Contents at a Glance

# Table of Contents

# Preface

The Python language has been around for a long time and worn many hats. Its original implementation was started by Guido van Rossum in 1989 as a tool for system administration as an alternative to Bash scripts and C programs.[1] Since its public release in 1991, it has evolved for use in a myriad of industries. These include everything from web-development, film, government, science, and business.[2]

I was first introduced to Python working in the film industry, where we used it to automate data management across departments and locations. In the last decade, Python has become a dominant tool in Data Science.

This dominance evolved due to two developments: the Jupyter notebook, and powerful third-party libraries. In 2001 Fernando Perez began the IPython project, an interactive Python environment inspired by Maple and Mathematica notebooks.[3] By 2014, the notebook-specific part of the project was split off as the Jupyter project. These notebooks have excelled for scientific and statistical work environments. In parallel with this development, third-party libraries for scientific and statistical computing were developed for Python. With so many applications, the functionality available to a Python programmer has grown immensely. With specialized packages for everything from opening web sockets to processing natural language text, there is more available than a beginning developer needs.

This project was the brainchild of Noah Gift.[4] In his work as an educator, he found that students of Data Science did not have a resource to learn just the parts of Python they needed. There were many general Python books and books about Data Science, but not resources for learning just the Python needed to get started in Data Science. That is what we have attempted to provide here. This book will not teach the Python needed to set up a web page or perform system administration. It is also not intended to teach you Data Science, but rather the Python needed to learn Data Science.

I hope you will find this guide a good companion in your quest to grow your Data Science knowledge.

# Example Code

Most of the code shown in examples in this book can be found on GitHub at: https://github.com/kbehrman/foundational-python-for-data-science.

---

1  https://docs.python.org/3/faq/general.html#why-was-python-created-in-the-first-place

2  https://www.python.org/success-stories/

3  http://blog.fperez.org/2012/01/ipython-notebook-historical.html

4  https://noahgift.com

# Figure Credits

| Figure | Credit Attribution |
|---|---|
| Cover | Boris Znaev/Shutterstock |
| Cover | Mark.G/Shutterstock |
| Figure 1-01 | Screenshot of Colab Dialogue © 2021 Google |
| Figure 1-02 | Screenshot of Renaming Notebook © 2021 Google |
| Figure 1-03 | Screenshot of Google Drive © 2021 Google |
| Figure 1-04 | Screenshot of Editing Text Cells © 2021 Google |
| Figure 1-05 | Screenshot of Formatting Text © 2021 Google |
| Figure 1-06 | Screenshot of Lists © 2021 Google |
| Figure 1-07 | Screenshot of Headings © 2021 Google |
| Figure 1-08 | Screenshot of Table of Contents © 2021 Google |
| Figure 1-09 | Screenshot of Hiding Cells © 2021 Google |
| Figure 1-10 | Screenshot of LaTeX Example © 2021 Google |
| Figure 1-11 | Screenshot of A Files © 2021 Google |
| Figure 1-12 | Screenshot of Upload Files © 2021 Google |
| Figure 1-13 | Screenshot of Mount Google Drive © 2021 Google |
| Figure 1-14 | Screenshot of Code Snippets © 2021 Google |

# Register Your Book

Register your copy of *Foundational Python for Data Science* on the InformIT site for convenient access to updates and/or corrections as they become available. To start the registration process, go to informit.com/register and log in or create an account. Enter the product ISBN **9780136624356** and click Submit. Look on the Registered Products tab for an Access Bonus Content link next to this product, and follow that link to access any available bonus materials. If you would like to be notified of exclusive offers on new editions and updates, please check the box to receive email from us.

# Acknowledgments

## About the Author

**Kennedy Behrman** is a veteran software engineer. He began using Python to manage digital assets in the visual effects industry and has used it extensively since. He has authored various books and training programs around Python education. He currently works as a senior data engineer at Envestnet.

*This page intentionally left blank*

# 3

# Sequences

*Errors using inadequate data are much less than those using no data at all.*

Charles Babbage

**In This Chapter**

- Shared sequence operations
- Lists and tuples
- Strings and string methods
- Ranges

In Chapter 2, "Fundamentals of Python," you learned about collections of types. This chapter introduces the group of built-in types called *sequences*. A sequence is an ordered, finite collection. You might think of a sequence as a shelf in a library, where each book on the shelf has a location and can be accessed easily if you know its place. The books are ordered, with each book (except those at the ends) having books before and after it. You can add books to the shelf, and you can remove them, and it is possible for the shelf to be empty. The built-in types that comprise a sequence are lists, tuples, strings, binary strings, and ranges. This chapter covers the shared characteristics and specifics of these types.

## Shared Operations

The sequences family shares quite a bit of functionality. Specifically, there are ways of using sequences that are applicable to most of the group members. There are operations that relate to sequences having a finite length, for accessing the items in a sequence, and for creating a new sequence based a sequence's content.

## Testing Membership

You can test whether an item is a member of a sequence by using the `in` operation. This operation returns `True` if the sequence contains an item that evaluates as equal to the item in question, and it returns `False` otherwise. The following are examples of using `in` with different sequence types:

```
'first' in ['first', 'second', 'third']
True

23 in (23,)
True

'b' in 'cat'
False

b'a' in b'ieojjza'
True
```

You can use the keyword `not` in conjunction with `in` to check whether something is absent from a sequence:

```
'b' not in 'cat'
True
```

The two places you are most likely to use `in` and `not in` are in an interactive session to explore data and as part of an `if` statement (see Chapter 5, "Execution Control").

## Indexing

Because a sequence is an ordered series of items, you can access an item in a sequence by using its position, or *index*. Indexes start at zero and go up to one less than the number of items. In an eight-item sequence, for example, the first item has an index of zero, and the last item an index of seven.

To access an item by using its index, you use square brackets around the index number. The following example defines a string and accesses its first and last substrings using their index numbers:

```
name = "Ignatius"
name[0]
'I'

name[4]
't'
```

You can also index counting back from the end of a sequence by using negative index numbers:

```
name[-1]
's'

name[-2]
'u'
```

## Slicing

You can use indexes to create new sequences that represent subsequences of the original. In square brackets, supply the beginning and ending index numbers of the subsequence separated by a colon, and a new sequence is returned:

```
name = "Ignatius"
name[2:5]
'nat'
```

The subsequence that is returned contains items starting from the first index and up to, but not including, the ending index. If you leave out the beginning index, the subsequence starts at the beginning of the parent sequence; if you leave out the end index, the subsequence goes to the end of the sequence:

```
name[:5]
'Ignat'
```

```
name[4:]
'tius'
```

You can use negative index numbers to create slices counting from the end of a sequence. This example shows how to grab the last three letters of a string:

```
name[-3:]
'ius'
```

If you want a slice to skip items, you can provide a third argument that indicates what to count by. So, if you have a list sequence of integers, as shown earlier, you can create a slice just by using the starting and ending index numbers:

```
scores = [0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18]
scores[3:15]
[3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14]
```

But you can also indicate the step to take, such as counting by threes:

```
scores[3:15:3]
[3, 6, 9, 12]
```

To count backward, you use a negative step:

```
scores[18:0:-4]
[18, 14, 10, 6, 2]
```

## Interrogation

You can perform shared operations on sequences to glean information about them. Because a sequence is finite, it has a length, which you can find by using the len function:

```
name = "Ignatius"
len(name)
8
```

You can use the `min` and `max` functions to find the minimum and maximum items, respectively:

```
scores = [0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18]
min(scores)
0
```

```
max(name)
'u'
```

These methods assume that the contents of a sequence can be compared in a way that implies an ordering. For sequence types that allow for mixed item types, an error occurs if the contents cannot be compared:

```
max(['Free', 2, 'b'])
--------------------------------------------------------------------
TypeError                            Traceback (most recent call last)
<ipython-input-15-d8babe38f9d9> in <module>()
----> 1 max(['Free', 2, 'b'])
TypeError: '>' not supported between instances of 'int' and 'str'
```

You can find out how many times an item appears in a sequence by using the `count` method:

```
name.count('a')
1
```

You can get the index of an item in a sequence by using the `index` method:

```
name.index('s')
7
```

You can use the result of the `index` method to create a slice up to an item, such as a letter in a string:

```
name[:name.index('u')]
'Ignati'
```

## Math Operations

You can perform addition and multiplication with sequences of the same type. When you do, you conduct these operations on the sequence, not on its contents. So, for example, adding the list [1] to the list [2] will produce the list [1,2], not [3]. Here is an example of using the plus (+) operator to create a new string from three separate strings:

```
"prefix" + "-" + "postfix"
'prefix-postfix'
```

The multiplication (*) operator works by performing multiple additions on the whole sequence, not on its contents:

```
[0,2] * 4
[0, 2, 0, 2, 0, 2, 0, 2]
```

This is a useful way of setting up a sequence with default values. For example, say that you want to track scores for a set number of participants in a list. You can initialize that list so that it has an initial score for each participant by using multiplication:

```
num_participants = 10
scores = [0] * num_participants
scores
[0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
```

# Lists and Tuples

Lists and tuples are sequences that can hold objects of any type. Their contents can be of mixed types, so you can have strings, integers, instances, floats, and anything else in the same list. The items in lists and tuples are separated by commas. The items in a list are enclosed in square brackets, and the items in a tuple are enclosed in parentheses. The main difference between lists and tuples is that lists are mutable, and tuples are immutable. This means that you can change the contents of a list, but once a tuple is created, it cannot be changed. If you want to change the contents of a tuple, you need to make a new one based on the content of the current one. Because of the mutability difference, lists have more functionality than tuples—and they also use more memory.

## Creating Lists and Tuples

You create a list by using the list constructor, `list()`, or by just using the square bracket syntax. To create a list with initial values, for example, simply supply the values in brackets:

```
some_list = [1,2,3]
some_list
[1, 2, 3]
```

You can create tuples by using the tuple constructor, `tuple()`, or using parentheses. If you want to create a tuple with a single item, you must follow that item with a comma, or Python will interpret the parentheses not as indicating a tuple but as indicating a logical grouping. You can also create a tuple without parentheses by just putting a comma after an item. Listing 3.1 provides examples of tuple creation.

Listing 3.1   **Creating Tuples**

```
tup = (1,2)
tup
(1,2)

tup = (1,)
tup
(1,)

tup = 1,2,
tup
(1,2)
```

> **Warning**
>
> A common but subtle bug occurs when you leave a trailing comma behind an argument to
> a function. It turns the argument into a tuple containing the original argument. So the second
> argument to the function `my_function(1, 2,)` will be `(2,)` and not 2.

You can also use the list or tuple constructors with a sequence as an argument. The following
example uses a string and creates a list of the items the string contains:

```
name = "Ignatius"
letters = list(name)
letters
['I', 'g', 'n', 'a', 't', 'i', 'u', 's']
```

## Adding and Removing List Items

You can add items to a list and remove items from a list. To conceptualize how it works, think of a
list as a stack of books. The most efficient way to add items to a list is to use the append method,
which adds an item to the end of the list, much as you could easily add a book to the top of a
stack. To add an item to a different position in the list, you can use the `insert` method, with the
index number where you wish to position the new item as an argument. This is less efficient than
using the append method as the other items in the list may need to move to make room for the
new item; however, this is typically an issue only in very large lists. Listing 3.2 shows examples of
appending and inserting.

Listing 3.2  **Appending and Inserting List Items**

```
flavours = ['Chocolate', 'Vanilla']
flavours
['Chocolate', 'Vanilla']

flavours.append('SuperFudgeNutPretzelTwist')
flavours
['Chocolate', 'Vanilla', 'SuperFudgeNutPretzelTwist']

flavours.insert(0,"sourMash")
flavours
['sourMash', 'Chocolate', 'Vanilla', 'SuperFudgeNutPretzelTwist']
```

To remove an item from a list, you use the pop method. With no argument, this method removes
the last item. By using an optional index argument, you can specify a specific item. In either case,
the item is removed from the list and returned.

The following example pops the last item off the list and then pops off the item at index 0. You can see
that both items are returned when they are popped and that they are then gone from the list:

```
flavours.pop()
'SuperFudgeNutPretzelTwist'
```

```
flavours.pop(0)
'sourMash'
```

```
flavours
 ['Chocolate', 'Vanilla']
```

To add the contents of one list to another, you use the `extend` method:

```
deserts = ['Cookies', 'Water Melon']
desserts
['Cookies', 'Water Melon']
```

```
desserts.extend(flavours)
desserts
['Cookies', 'Water Melon', 'Chocolate', 'Vanilla']
```

This method modifies the first list so that it now has the contents of the second list appended to its contents.

### Nested List Initialization

There is a tricky bug that bites beginning Python developers. It involves combining list mutability with the nature of multiplying sequences. If you want to initialize a list containing four sublists, you might try multiplying a single list in a list like this:

```
lists = [[]] * 4
lists
[[], [], [], []]
```

This appears to have worked, until you modify one of the sublists:

```
lists[-1].append(4)
lists
[[4], [4], [4], [4]]
```

All of the sublists are modified! This is because the multiplication only initializes one list and references it four times. The references look independent until you try modifying one. The solution to this is to use a list comprehension (discussed further in Chapter 13, "Functional Programming"):

```
lists = [[] for _ in range(4)]
lists[-1].append(4)
lists
 [[], [], [], [4]]
```

## Unpacking

You can assign values to multiple variables from a list or tuple in one line:

```
a, b, c = (1,3,4)
a
1
```

```
b
3

c
4
```

Or, if you want to assign multiple values to one variable while assigning single ones to the others, you can use a * next to the variable that will take multiple values. Then that variable will absorb all the items not assigned to other variables:

```
*first, middle, last = ['horse', 'carrot', 'swan', 'burrito', 'fly']
first
['horse', 'carrot', 'swan']

last
'fly'

middle
'burrito'
```

## Sorting Lists

For lists you can use built-in sort and reverse methods that can change the order of the contents. Much like the sequence min and max functions, these methods work only if the contents are comparable, as shown in these examples:

```
name = "Ignatius"
letters = list(name)
letters
['I', 'g', 'n', 'a', 't', 'i', 'u', 's']

letters.sort()
letters
['I', 'a', 'g', 'i', 'n', 's', 't', 'u']

letters.reverse()
letters
['u', 't', 's', 'n', 'i', 'g', 'a', 'I']
```

# Strings

A string is a sequence of characters. In Python, strings are Unicode by default, and any Unicode character can be part of a string. Strings are represented as characters surrounded by quotation marks. Single or double quotations both work, and strings made with them are equal:

```
'Here is a string'
'Here is a string'

"Here is a string" == 'Here is a string'
True
```

If you want to include quotation marks around a word or words within a string, you need to use one type of quotation marks—single or double—to enclose that word or words and use the other type of quotation marks to enclose the whole string. The following example shows the word *is* enclosed in double quotation marks and the whole string enclosed in single quotation marks:

```
'Here "is" a string'
'Here "is" a string'
```

You enclose multiple-line strings in three sets of double quotation marks as shown in the following example:

```
a_very_large_phrase = """
Wikipedia is hosted by the Wikimedia Foundation,
a non-profit organization that also hosts a range of other projects.
"""
```

With Python strings you can use special characters, each preceded by a backslash. The special characters include \t for tab, \r for carriage return, and \n for newline. These characters are interpreted with special meaning during printing. While these characters are generally useful, they can be inconvenient if you are representing a Windows path:

```
windows_path = "c:\row\the\boat\now"
print(windows_path)


ow heoat
     ow
```

For such situations, you can use Python's raw string type, which interprets all characters literally. You signify the raw string type by prefixing the string with an r:

```
windows_path = r"c:\row\the\boat\now"
print(windows_path)
c:\row\the\boat\now
```

As demonstrated in Listing 3.3, there are a number of string helper functions that enable you to deal with different capitalizations.

Listing 3.3   **String Helper Functions**

```
captain = "Patrick Tayluer"
captain
'Patrick Tayluer'

captain.capitalize()
'Patrick tayluer'

captain.lower()
'patrick tayluer'

captain.upper()
'PATRICK TAYLUER'
```

```
captain.swapcase()
'pATRICK tAYLUER'

captain = 'patrick tayluer'
captain.title()
'Patrick Tayluer'
```

Python 3.6 introduced format strings, or f-strings. You can insert values into f-strings at runtime by using replacement fields, which are delimited by curly braces. You can insert any expression, including variables, into the replacement field. An f-string is prefixed with either an F or an f, as shown in this example:

```
strings_count = 5
frets_count = 24
f"Noam Pikelny's banjo has {strings_count} strings and {frets_count} frets"
'Noam Pikelny's banjo has 5 strings and 24 frets'
```

This example shows how to insert a mathematic expression into the replacement field:

```
a = 12
b = 32
f"{a} times {b} equals {a*b}"
'12 times 32 equals 384'
```

This example shows how to insert items from a list into the replacement field:

```
players = ["Tony Trischka", "Bill Evans", "Alan Munde"]
f"Performances will be held by {players[1]}, {players[0]}, and {players[2]}"
'Performances will be held by Bill Evans, Tony Trischka, and Alan Munde'
```

## Ranges

Using range objects is an efficient way to represent a series of numbers, ordered by value. They are largely used for specifying the number of times a loop should run. Chapter 5 introduces loops. Range objects can take start (optional), end, and step (optional) arguments. Much as with slicing, the start is included in the range, and the end is not. Also as with slicing, you can use negative steps to count down. Ranges calculate numbers as you request them, and so they don't need to store more memory for large ranges. Listing 3.4 demonstrates how to create ranges with and without the optional arguments. This listing makes lists from the ranges so that you can see the full contents that the range would supply.

Listing 3.4  **Creating Ranges**

```
range(10)
range(0, 10)

list(range(1, 10))
[1, 2, 3, 4, 5, 6, 7, 8, 9]
```

```
list(range(0,10,2))
[0, 2, 4, 6, 8]

list(range(10, 0, -2))
[10, 8, 6, 4, 2]
```

## Summary

This chapter covers the import group of types known as sequences. A sequence is an ordered, finite collection of items. Lists and tuples can contain mixed types. Lists can be modified after creation, but tuples cannot. Strings are sequences of text. Range objects are used to describe ranges of numbers. Lists, strings, and ranges are among the most commonly used types in Python.

## Questions

1. How would you test whether `a` is in the list `my_list`?

2. How would you find out how many times `b` appears in a string named `my_string`?

3. How would you add `a` to the end of the list `my_list`?

4. Are the strings `'superior'` and `"superior"` equal?

5. How would you make a range going from 3 to 13?

# Index