


Python for Data Analysis

Lecture Notes

By: [Alexander Clark](#)¹
Columbia University SPS

This version: May 8, 2022

¹ac4725@columbia.edu 

Contents

Preface	vii
I Programming	1
1 Input, Processing, and Output	3
1.1 The Hello World Program	3
1.2 Variables	3
1.3 Comments	4
1.4 Data Types and Conversion	5
1.5 Input	5
1.6 Calculations	5
2 If, Elif, Else	7
2.1 If Statements	7
2.2 If Elif Else Statements	7
2.3 Logical Operators	8
2.4 Boolean Variables	9
3 Lists	11
3.1 Defining Lists	11
3.1.1 Indexing, Slicing, Mutating	11
4 Repetition Structures	15
4.1 While Loops	15
4.2 For Loops	16
4.2.1 Nested Loops	17
5 Functions	19
5.1 Detour: Callables and Functions	19
5.2 Function Basics	19
5.2.1 Naming	20
5.3 Arguments and Parameters	20
5.3.1 A mutability gotcha	21
5.4 Annotation and Documentation	21
5.5 Local and Global Variables and Namespaces	22
5.6 Anonymous Functions	22
5.7 Higher-Order Functions Like <code>map</code>	23
5.8 Examples	23
5.8.1 Exercise!	23
6 Exceptions	25
6.1 Exceptions	25

7	More Data Structures	29
7.1	Tuples	29
7.1.1	Tuple Assignment	29
7.1.2	Unpacking	30
7.1.3	Tuples vs. Lists	30
7.2	The <code>in</code> Operator	30
7.3	List Methods	31
7.4	Strings	32
7.5	Dictionaries	33
7.6	Sets	34
8	Modules	35
8.1	Storing Functions in Modules	35
8.2	Using the <code>random</code> module and <code>matplotlib</code>	36
9	Object Oriented Programming	37
9.1	Object Oriented Programming	37
9.1.1	Classes	37
9.2	Inheritance	39
9.2.1	Polymorphism (G§11.2)	39
10	IPython	41
10.1	IPython and Jupyter	41
10.1.1	IPython Basics	41
10.1.2	Extensions	41
10.1.3	Peak Under the Hood: Grading	41
II	Data	43
11	NumPy	45
11.1	The Array	45
11.2	Indexing	45
11.2.1	Functions	46
11.2.2	Linear Algebra	46
12	Pandas: Series and DataFrames	49
12.1	Series	49
12.1.1	Series Functionality I	50
12.2	DataFrames	51
12.2.1	Indexing	52
12.3	Summarizing and Computing Descriptive Stats	52
12.4	Applications	52
12.4.1	Interview Question	52
13	Pandas: Join, Merge, and Other Manipulation	55
13.1	Application: Primitive Pandas	55
13.2	The Basic Join	55
13.3	Data Aggregation and Group Operations	56
13.3.1	GroupBy	56
13.4	Concat and Append	57
13.4.1	Application: What precedes sleeplessness?	58
13.5	Merge	58
13.6	Pivot Tables	60
13.6.1	Crosstabs	60

14 Data Visualization	61
14.1 MATLAB Interface	61
14.1.1 Saving Figures	62
14.1.2 Special Plots: Scatter, Histogram, and Bar Plots	62
14.1.3 Customizations	62
14.1.4 Pandas Integration	62
14.2 Object-oriented Interface	63
14.2.1 Introduction to the OO Interface	63
14.2.2 Subplots with <code>plt.subplots()</code>	67
15 Time and Dates	69
15.1 Using the <code>datetime</code> module	69
15.1.1 Dates and Datetimes	69
15.1.2 Unix Timestamps	69
15.1.3 Time Deltas and Relative Deltas	70
15.1.4 Date Strings	70
15.1.5 Pandas and Numpy	71
15.2 Analysis	71
15.2.1 Application 1	72
15.2.2 Application 2	72
16 Python for Excel	73
16.1 Pandas	73
16.1.1 Reading	73
16.1.2 Writing	74
16.2 OpenPyXL	74
16.2.1 Reading	74
16.2.2 Writing	75
16.2.3 Styles	76
16.2.4 Inserting Charts	80
17 Applications	83
17.1 Wild Data Redux	83
17.1.1 Missing Data	84
17.1.2 Categorical	86
17.2 Inference and Experiments	87
17.2.1 Motivation/Soapbox	87
17.2.2 Experiments	89

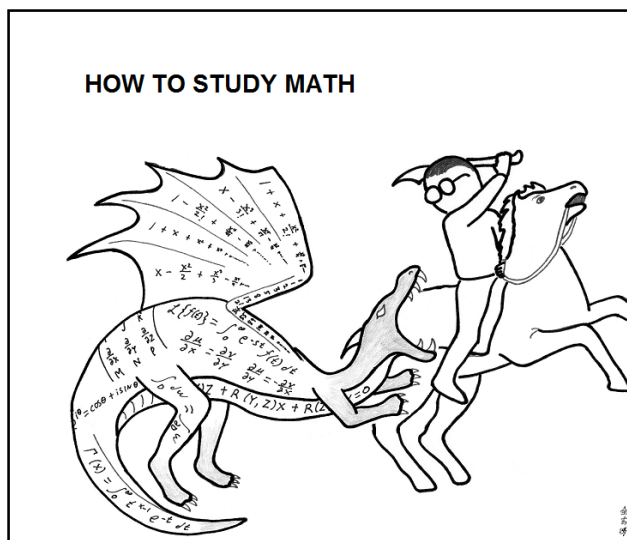
Preface

These are my lecture notes for Python for Data Analysis, which I consider a work in progress. They're not written to replace class attendance, so you might not find them self-contained.

The main references for these notes are Lubanovic 2019 and VanderPlas 2016a. There are many other fantastic books and resources for self-studying Python. I encourage you to find as many supplementary sources as you find helpful, but to return to these notes to help focus on what will be important for doing well on assignments and exams in this course.

I've seen this material presented in many different orders. What seems unavoidable in any bearable introduction to Python is that some concept might be slipped in informally before we regroup for a formal explanation of how some data type works or why we keep using words like *object*, *method*, or *callable*. This isn't so different than learning a language through immersion. That's to say we'll be immersed and you might find yourself flipping backwards or forwards in these notes. You could also compare it to bulking and cutting in bodybuilding. Let me know if you see anything that could use restructuring.

Famous mathematician Paul Halmos counseled students not to study math passively. With Python, this is also good advice. Ease into the language, but don't remain satisfied with running others' code or making only minimal edits. Work from scratch and relish the detours.



Don't just read it; fight it!

Source: AbstruseGoose.com --- Paul R. Halmos

Part I

Programming

Chapter 1

Input, Processing, and Output

1.1 The Hello World Program

To get started in any language, printing “Hello, World!” might be the first step.¹

In Python, we can print an input using the `print()` function. We simply pass our desired input within the parentheses, and Python will print the value.

We can enter text as a **string**. Text entered inside single, double, or triple quotations is interpreted as a string.

```
1 print('Hello, World')
2 print("Hello, World!")
3 print("""Hello,
4 World!""")
```

In Python 2, the syntax would have been `print 'Hello, World!'` without the parentheses. You shouldn't use Python 2 or bother learning its variations, but this is a good difference to understand. If you see this old syntax in an old Stack Exchange post or anywhere else, let that be a tipoff that you might be looking at Python 2 code, which might behave differently.

1.2 Variables

A **variable** holds a value. It can be a string, a number, or perhaps a more complicated data type. Variable assignment is done with the equals sign, `=`.

```
1 greeting = "Hello, World!"
2 my_favorite_number = 91
```

Now compare the output you get from the following.

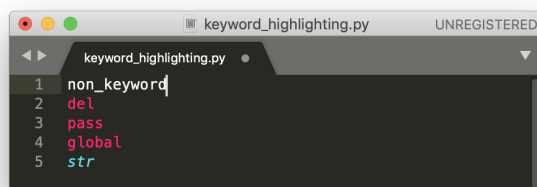
```
1 print(greeting)
2 print('greeting')
3 print("Hello, World!")
4 print(91)
5 print(my_favorite_number)
6 print('my_favorite_number')
7 print("My favorite number is", my_favorite_number)
8 print("My favorite number is ", my_favorite_number)
```

PEP 8 addresses variable names [here](#). Use lowercase and underscores. This is good advice, but I don't see any reason to be too wedded to this. If you want to assign a matrix to a variable, it's reasonable to use an uppercase letter as the variable name. There, a math convention overrides a Python convention.

More importantly, avoid Python key words in your variable names. [Here](#) is a list of key words which have a specific meaning in code. See Lubanovic 2019 page 27 for a complete list of rules for variable names. An

¹See [the Wikipedia page](#) for Hello, World.

IDE or code editor will make this easier for you by highlighting keywords that you shouldn't use for variable names.



While you simply can't assign any value to `class`, you can assign a value to `str`. Even though it's permitted, I would suggest you never assign a value to `str` because it can be used to confirm something is a string type. The same goes for keywords like `int`. Consider the program below.

```
1 s = 'Python'
2 is_string = type(s) == str
```

The above works so that `is_string` holds a value `True` if and only if `s` is a string. In this case, `is_string = True`. The program below also runs, but it does something different. `str` holds the value 10 so the last line checks if the type of `s` is equal to ten. Now, `is_string` holds the value `False`.

```
1 s = 'Python'
2 str = 10
3 is_string = type(s) == str
```

Finally, consider how the variables chapter in Lubanovic 2019 opens with a quote from Proverbs 22:1. “A good name is rather to be chosen than great riches.” In its Biblical context, this expresses the idea that a good reputation is better than money and that the name is an expression of the inner character of its bearer. Lubanovic 2019 is stressing the importance of variable names in a different way. Well-named variables help the reader and user of your code. In a way, instead of indicating the character of the person, variable names actually shape the character of your code. Be thoughtful of this, so if you're using `n` for the number of apples, perhaps you can spare the keystrokes for `num_apples`. We might complete the Biblical analogy by saying that readable code with well-named variables is better than confusing code that might run faster or have been written faster. This concern for readability brings us to our next topic, comments.

1.3 Comments

Comments are briefly addressed in Lubanovic 2019 Chapter 4.

Commenting your code is helpful if you care about your colleagues or your future self. Comments should add clarity to the intention and workings of code. A comment is a piece of code that isn't actually executed—it's a comment left for the reader or the person who inherits and modifies your code. Everything after a `#` will be ignored by the Python interpreter.

```
1 # This will print a greeting.
2 print('Hello, World!')
```

You might also use end-line comments like the following

```
1 print('Hello, World!') # Prints a greeting
```

PEP 8 addresses comments [here](#). I don't intend to grade based on the stylistic orthodoxy of your comments, but spaces are free so I do recommend `# Comments like this` instead of `#Comments like this`.

Perfect comment technique does not correct for bad code though. Compare the following blocks of code.

```
1 x = 90 # Wins
2 y = 10 # Losses
3 z = x/y # Win Loss Ratio
4 a = 100 * x/(y+x) # Winning Percentage
```

```

1 wins = 90
2 losses = 10
3 win_loss_ratio = wins/losses
4 winning_percentage = 100 * wins / (losses + wins)

```

The first block is commented and the second is not. Still, the second code block is much better because the variable names are chosen so you don't *need* comments. Good naming becomes even more important as the program becomes longer and the variable is used over and over. See also the discussion in Long 2021.

1.4 Data Types and Conversion

To start, we are concerned with strings, integers, and floats. In Python, these are classes `str`, `int`, and `float`. You can check the type of variable or value using `type()`.

```

1 string_example = ''
2 int_example = -1
3 float_example = -1.

```

Some types can be converted by using `str()`, `int()`, or `float()`. You might informally call these functions, but they actually aren't because function has a special meaning that's reserved for other things. We'll call them callables, which is a more general category.

```

1 print(type(1))
2 print(type(str(1)))

```

1.5 Input

It's not that common for a data science workflow, but you can read input using `input()`. The input is always read in as a string.

```

1 favorite_color = input("What is your favorite color?")
2 favorite_number = input("What is your favorite number?")
3 attending_in_person = input("I am attending class in person.")

```

1.6 Calculations

You can use Python as a calculator. Below is the list of operations and symbols.

Operator	Description
+	Addition
-	Subtraction
*	Multiplication
/	Floating point (normal) division
//	Integer (truncating) division
%	Modulus (remainder)
**	Exponentiation

What might stand out is

- Exponentiation is done with `**`, not `^`.
- Integer division (rounds down) is done with `//`.
- The remainder of x divided by y can be found with `x % y`, which might be read as x modulo y .

If a float is involved in an operation, the result will also be a float, as in `type(2 + 2.)`.

Chapter 2

If, Elif, Else

Reference: Lubanovic Chapter 4

Decision structures allow a program to have more than one path of execution. The path depends on condition. The condition is either True or False, and so can be represented by a Boolean variable.

2.1 If Statements

Here's a joke. A programmer is going to the grocery store and his partner says, "Buy a gallon of milk, and if there are eggs, buy a dozen." The programmer comes home with 13 gallons of milk.

Or consider the logical inference if you ask, "Is it raining?" and get a reply, "Not hard."

```
1 if 2 + 2 > 4:
2     print("Pigs can fly.")
3
4 if 2 + 2 == 4:
5     print("Pigs cannot fly.")
6
7 if 'a' < 'b':
8     print("It is true that 'a' is less than 'b'.")
9
10 if 'a' < 'A':
11     print("It is true that 'a' is less than 'A'.")
12
13 if 'goon' == 'Goblin':
14     print("A goon is a goblin.")
```

If statements like the above rely on *relational operators* (see Table 3-1 Gaddis p.112).

```
1 if 2 == 2.:
2     print("The integer and the float are equal.")
3
4 if 2 is 2.:
5     print("The integer and the float are the same object in memory.")
```

2.2 If Elif Else Statements

The natural counterpart to `if` is `else`. The code under `else` simply executes when the if condition is not satisfied. We also have `elif` (or else if) to help in the intermediate case, where we want another block of code to be run when the if condition is not satisfied *and* some other condition is satisfied.

Compare this program with a simpler version using else if conditions.

```
1 x = 1
2 if x == 0:
3     print("zero")
4 else:
```

```

5     if x < 0:
6         print("negative")
7     else:
8         print("positive")

```

```

1 x = 1
2 if x == 0:
3     print("zero")
4 elif x < 0:
5     print('negative')
6 else:
7     print('positive')

```

Compare the output from the following programs.

```

1 num = 0
2
3 if num < 1:
4     print(num)
5     num = num + 2
6 if num > 0:
7     print(num, '!')
8     num = num - 1000
9 if True:
10    print(num, '?')

```

```

1 num = 0
2
3 if num < 1:
4     print(num)
5     num = num + 2
6 elif num > 0:
7     print(num, '!')
8     num = num - 1000
9 else:
10    print(num, '?')

```

2.3 Logical Operators

Suppose you want to execute some code if a number x is between 10 and 20. You could use *nested* if statements.

```

1 x = 14
2 if x >= 10:
3     if x <= 20:
4         print("x is between 10 and 20.")

```

But you might prefer to base your if statement off of one compound Boolean expression. For these, we need logical operations. They are

- Logical *and*: **and** (& for bitwise)
- Logical *or*: **or** (| for bitwise)
- Logical *negation*: **not** (~ bitwise)

Observe the following will give equivalent output.

```

1 not 1 > 2
2 not (1 > 2)
3 not(1 > 2)

```

Challenge: What do you expect from **not not** (1 **or** False)?

The usefulness of these logical operators is in chaining together several boolean expressions, reducing them to one.

Consider the following three blocks of code, which do the same thing.


```
1 if x >= 10:
2     if x <= 10:
3         print("it's 10!")

1 if (x >= 10) and (x <= 10):
2     print("it's 10!")

1 if 10 <= x <= 10:
2     print("it's 10!")
```

The latter two eliminate the need for nesting if statements. This helps readability. As the [Zen of Python](#) says, “flat is better than nested.”

2.4 Boolean Variables

Finally, a Boolean variable simply references a logical True or False.

```
1 # Option Value
2 market_value = 10
3 strike_price = 9
4 option_has_value = market_value > strike_price
5
6 if option_has_value:
7     print("We're in the money.")
8
9 # Check data type
10 print(type(option_has_value))
11 print(type(False))
```

Note the above program could be shortened. We could accomplish the same thing in fewer lines and with fewer characters. We could cut out the variable `option_has_value` and place the logical condition it represents directly in the if statement. If you’re playing [code golf](#), this is a good idea. But there’s a reason code golf is described as “recreational” programming. It’s a fun challenge, but it’s a concept independent of readable code. Define extra variables to simplify your own programs and make them more readable!

Chapter 3

Lists

Reference: Lubanovic 2019 Chapter 7

In this chapter we introduce lists and this occasions the discussion of a new idea, (im)mutability. We cover just the basics to help prepare us for for loops.

3.1 Defining Lists

We have worked with strings, integers, floats, and booleans so far. Now, we introduce a compound data type the list. Lists are also a sequence object, meaning they are ordered.

A list is constructed by separating one or more objects with commas and placing them in square brackets. For example, we can define a list as

```
1 # Our first list
2 floor = ['yoga', 'meditation', 'stretching', 'strength', 'cardiovascular']

1 # More lists
2 aquatic = []
3 cycling = ['cycling']
```

Above, `[]` created an empty list. You could also use, and you might prefer `list()`. The reason you might prefer `list()` is because explicit is better than implicit¹, and this makes it explicit that you're creating a list.

```
1 # More lists
2 treadmill_based = ['running', 'walking', 'bootcamp']
```

A useful quality of lists is that they can be of mixed type.

```
1 # More lists
2 fine_list = [0, 'milk', cycling]
```

3.1.1 Indexing, Slicing, Mutating

You can obtain the element at a certain place, or *index*, in a list by suffixing the list with that index number inside square brackets. Python uses zero-based indexing. So the n^{th} item is at index $n - 1$.

```
1 # Get specific items
2 fine_list = [0, 'milk', cycling]
3
4 # Which of these print statements will raise an error?
5 print(fine_list[0])
6 print(fine_list[1])
7 print(fine_list[2])
8 print(fine_list[3])
9 print(fine_list[-1])
```

¹From the Zen of Python.

We can also count from the end of the list to find an item, as hinted by the `fine_list[-1]` above. The index `-1` will identify the last element. So, in some sense, you might think that negative indexing is not zero-based. If that's confusing, just recall $-0 = 0$ and so `fine_list[-0]` is actually the same as `fine_list[0]`.

The following graphic helps explain how lists (and strings) are indexed.

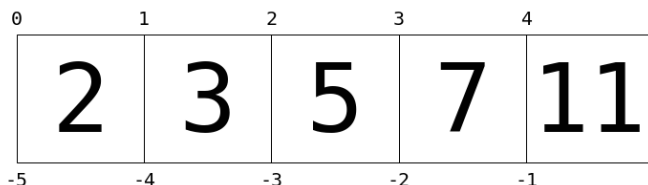


Figure 3.1: Indexing (From VanderPlas 2016b).

While indexing pulled out a single element, we can also *slice* to pull out a sublist.

```
1 # Get specific items
2 ['Jack', 'Jill', 'and', 'hill', 'over', 'ran', 'the']
3 print(fine_list[0:2])
4 print(fine_list[-2:])
```

The first print statement will print elements 0 and 1 in a list. The last print statement will print a list of just the last two elements. Why is slicing done this way? Why does the numbering start at zero and why is the upper bound excluded? It's more of a computer science convention than a Python-specific one. Ramalho 2015 (freely downloadable from Columbia Library) addresses this on page 33 and Dijkstra provides some reasoning [here](#). Those reasons are basically:

1. The length of the list is more obvious. `fine_list[0:2]` contains $2-0=2$ elements. And `fine_list[:4]` contains four elements.
2. It's easy to partition a list. `fine_list[:3]` and `fine_list[3:]` don't overlap.

Lists are *mutable*, meaning that we can change their contents.

```
1 # Friend drama
2 my_best_friends = ['Richard', 'Spiro']
3 print(my_best_friends)
4 my_best_friends[1] = 'Gerald'
5 print(my_best_friends)
```

Mutability has consequences when you are assigning one variable to be equal another. By changing `b`, which is made from `a`, we also change `a`.

```
1 a = [1,3,5]
2
3 b = a
4 b[2] = 99
5
6 print(a)
7 print(b)
```

This was *not* the case for immutable objects like integers, floats, and strings.

```
1 a = 1
2 b = a
3 a += 1
4
5 # b and a are not the same
6 print(a)
7 print(b)
```

However, we don't see the same effect on mutable objects when reassigning them.

```
1 a = [1,3,5]
2
3 b = a
4 b = b[0:1]
5
6 # b and a are no longer pointing to the same list
7 print(a)
8 print(b)
```

Just when “mutating” them.

```
1 a = [1,3,5]
2
3 b = a
4 b = b.append(7)
5
6 # b and a are pointing to the same list
7 print(a)
8 print(b)
9
10
11 # This also does the same
12 a = [1]
13 b = a
14
15 a += [1]
16
17 # b and a are pointing to the same list
18 print(a)
19 print(b)
```


Chapter 4

Repetition Structures

Repetition structures (loops) are one of the best justifications for moving from Excel to Python (though they are not unique to Python). A “program” in Excel is a series of keystrokes and clicks. You might create a report for your company that is specific to one market and you’ll need to replicate the same report for a different market. Perhaps you could write a macro, but I think you’ll find working in Python to be easier. In Python, we can do this in a loop. We can have a program that makes the report and we can iterate through the different markets to apply the program to each and create the specific reports (using a for loop).

We also introduce the idea of an *iterable* object.

4.1 While Loops

The **while** loop is a condition-controlled loop. Figure 4.1 illustrates the logic well.

The statement inside a while loop executes as long as the condition evaluates to True.

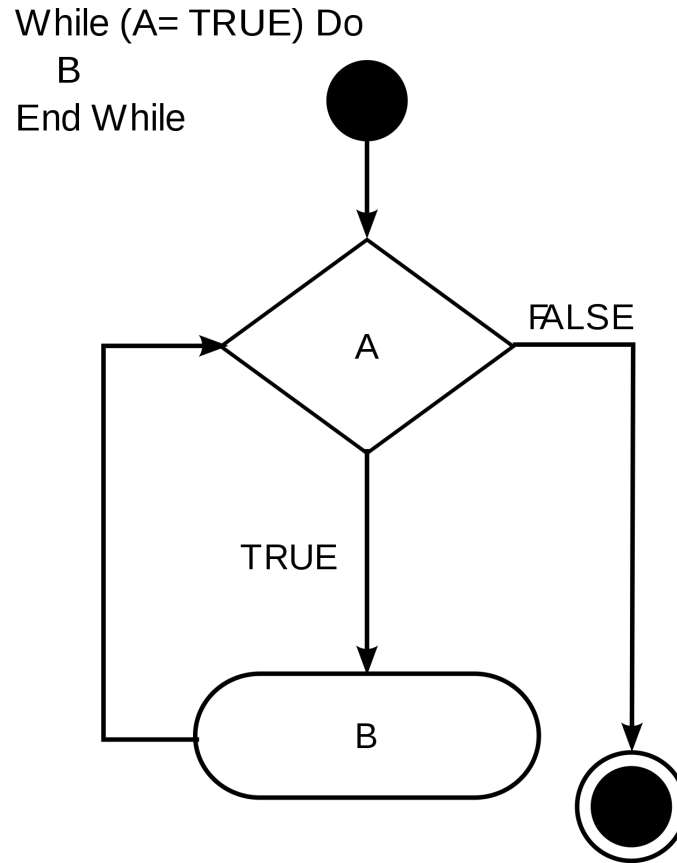
```
1 # Rest on the seventh day
2 day_of_week = 1
3 while day_of_week < 7:
4     print('work')
5     day_of_week += 1
```

Beware the infinite loop. Be confident that your test condition will have a way of becoming False, otherwise you might notice that your program never finishes.

Let’s consider an infinite series, $\sum_{i=0}^{\infty} 2^{-i} = 1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \dots = 2$. We could try to calculate this with a while loop if we didn’t know the sum converged to two.

```
1 # Sum of a geometric series
2 the_sum = 0
3 idx = 0
4 increment = 2 ** -idx
5
6 while increment > 0:
7     # Increase the sum by the current increment
8     the_sum += increment
9
10    # Advance the index in the sum and calculate a new increment
11    idx += 1
12    increment = 2 ** (-idx)
```

Does this loop make you nervous? In fact, the loop will terminate because we will eventually hit machine zero. I found the loop to terminate at the increment 2^{-1075} and the resulting sum to be two. But it should make you nervous. You might instead decide on some level of precision and use a test condition like `increment > .0005`, in which case you could find a bound on the error with some math. Doing that math is not part of this course, (un)fortunately.

Figure 4.1: While loop logic ([Wikipedia](#)).

4.2 For Loops

For loops execute the attached code based on some iteration. The code might depend on a variable that is actually changing with the iteration.

```

1 #Dr. Seuss
2 for item in [1,2,'red','blue']:
3     print(item, 'fish')
  
```

The general structure is

```

1 for <iteration_variable> in <iterable>:
2     <program>
  
```

The iteration variable should be chosen to maximize readability.

Sometimes you find definitions for **iterable** that are circular, like “objects that can be iterated over.” It’s something where we can traverse through distinct values. So far, the iterables we have seen are strings and lists. A for loop will iterate over a list one element at a time, according to the ordering in the list. A for loop iterates over a string one character at a time. There are several more types of iterable objects, but for now the only additional one we’ll introduce is that created by `range(start, stop, step)`, which allows you to iterate over integers from `start` (inclusive) to `stop` (exclusive) according to the `step` parameter. By default, `start = 0` and `step = 1` so these are optional arguments. You can check what is included in a particular call to `range()` by examining `list(range(start, stop, step))`.

The for clause tells Python to execute the statement once for every item in the iterable, which is in this case the object `[1,2,'red','blue']`. This is a *list*, a kind of compound data object that can store other objects in

sequence by separating them with commas and putting them inside square brackets. Below we use `range(5)` instead of a list, which you can think of as generating an iterable object of integers from 0 to 4 (5 is not included, but the length is 5).

Or you might just want to execute a specific set of statements some number of times.

```
1 # Tubthumping by Chumbawamba
2 for item in range(5):
3     print("I get knocked down, but I get up again.")
4     print("You are never gonna keep me down.")
```

For both of the examples above, `item` is the *variable*. However, notice the print statement only depends on the variable in the first example.

You can even iterate over the characters in a string.

```
1 # Cheer
2 word = 'PYTHON'
3 for char in word:
4     print("Give me a",char, '!')
5     print("    ", char)
6 print("What's that spell?")
7 print("    ",word)
```

4.2.1 Nested Loops

A loop that is inside another is called *nested*.

Think about how Python executes code line by line. What order of output do you expect from this program?

```
1 # Nested Loop
2 for x in ['wee', 'bee']:
3     for y in ['bop', 'dop']:
4         print(x,y)
```

If it helps, the above is a simplification of the following.

```
1 # Nested Loop
2 x = 'wee'
3 for y in ['bop', 'dop']:
4     print(x,y)
5
6 x = 'bee'
7 for y in ['bop', 'dop']:
8     print(x,y)
```


Chapter 5

Functions

Reference: [Lubanovic 2019 Chapter 9](#)

Functions make your code easier to read and write by performing a certain task based on some number of arguments the function takes. Whenever you find yourself copying and pasting code in a program, you should ask yourself if you shouldn't be using a function instead.

5.1 Detour: Callables and Functions

It's a common mistake to call more things functions than actually are functions. There are other things like methods and classes you'll be formally introduced to later and that have a similar feel. Classes and functions are both *callable*s, just meaning something you can call by using parentheses and sometimes arguments inside those parentheses.

Consider the code and output below (pasted from the Terminal in VSCode).

```
1 >>> list
2 <class 'list'>
3 >>> len
4 <built-in function len>
5 >>> callable(len)
6 True
7 >>> callable(list)
8 True
9 >>> callable
10 <built-in function callable>
11 >>> callable(callable)
12 True
```

5.2 Function Basics

In math, you might define a function $f(x) = x^2$. Just as f does something with x , a function like `print` does something with whatever input we pass inside the parentheses. However, an argument/parameter is not always necessary or even allowed depending on the function. In the example below, we define our first function.

```
1 def your_annoying_friend():
2     print("crypto")
```

A function definition is begun with the keyword `def`. Then you give the function name and input arguments inside parentheses. After the parentheses there is a colon. Then the code dictating what the function does is indented below.

Our function `your_annoying_friend` takes no arguments. No matter how many times you call `your_annoying_friend`, all it does is print `'crypto'`. Try passing an argument inside, like `your_annoying_friend('new conversation topic')`, but you'll only get an error because the function is not defined to accept arguments.

Further this function actually *returns* nothing because there is no *return statement*. Try running `a = your_annoying_friend()`. The variable `a` has the value `None`, the `NoneType` (see Lubanovic 2019 page 144 for a brief discussion).

We'll introduce return statements shortly, but first, let's include an argument in a new function. Consider the cheer we made in Section 4.2. We can make this into a function so that it works for any value of `word`.

```
1 # Cheer
2 word = 'PYTHON'
3 for char in word:
4     print("Give me a",char, '!')
5     print("    ", char)
6 print("What's that spell?")
7 print("    ",word)

1 # Cheer Function
2 def cheer(word):
3     ''' Doc string '''
4     for char in word:
5         print("Give me a",char, '!')
6         print("    ", char)
7     print("What's that spell?")
8     print("    ",word)
```

A function will return a value once it reaches a return statement, begun with the keyword `return`. Once the return statement is reached, the function quits executing and any leftover bits of the program are abandoned. Call the function one below. You'll find that `"testing"` is never printed because it comes after `return 1`. We can say this function accepts no parameters and always does two things: *prints "Coming right up."* and *returns 1*.

```
1 def one():
2     print("Coming right up.")
3     return 1
4     print("testing")
```

5.2.1 Naming

The same rules for naming variables apply to naming functions. Again, see [PEP 8](#). Readability counts.

5.3 Arguments and Parameters

I might slip, but there is technically a difference between *arguments* and *parameters*. Lubanovic 2019 tells us the values you pass into the function are the arguments. Those values are copied the corresponding *parameters* inside the function. In the below, `add` is defined and called with arguments 1 and 2, and these are copied to the parameters `x` and `y`.

```
1 def add(x,y):
2     return x + y
3
4 add(1,2)
```

The order of your arguments matters.

```
1 def investment_strategy(buy, sell):
2     print("buy", buy)
3     print("sell", sell)
4
```

Calling `investment_strategy('high', 'low')` and `investment_strategy('low', 'high')` produce very different ideas. These arguments are *positional arguments*, in that they are copied to their parameters based on the ordering. If that produces confusion, you can instead use *keyword arguments*.

```
1 investment_strategy(sell = 'high', buy = 'low')
2
```

You can also specify *default* parameter values so that you don't have to pass an argument. This is demonstrated in the definition below.

```
1 def my_favorite_class(a = 'Python'):
2     print(a, "!", sep = '')
```

5.3.1 A mutability gotcha

The default arguments are evaluated once when a function is defined.¹ This can create something unexpected when using a mutable default argument—if you mutate the argument, it will stay mutated for future use of the function.

```
1 def risky_function(x, l = []):
2     l.append(x)
3     return l
```

Run `risky_function('a')` twice. You might expect `['a']` to be returned in both instances, but you will get `['a', 'a']` on the second call. You can avoid this behavior with something like the following.

```
1 def workaround(x, l = None):
2     if l == None:
3         l = list()
4     l.append(x)
5     return l
```

5.4 Annotation and Documentation

A function can include a documentation string (docstring) and annotation regarding the intended argument data types. Annotations for a parameter are to be preceded with a colon. An annotation for the output can be created with `->` before the annotation, between the closed parenthesis and colon. Python simply stores this information. There are no checks or enforcement. Below, the default argument for `y` doesn't even obey the annotation.

```
1 def annotated_function(x:str, y:'int > 0' = -1) -> bool:
2     "Sample Docstring"
3     return x
4
5 # Violates the annotation but still runs
6 annotated_function(1, -99)
```

If you want checks enforced, you will have to code them in yourself. One way to do that is by using an `assert` statement inside your function. With `assert` and then a boolean expression, your code will break if the expression is false.

```
1 def assertive_function(x: int):
2     assert type(x) == int
3     return x
```

This doesn't help the user of your code very much. So instead, you might want to create an exception of your own to explain what went wrong. This is done with `raise`, which will be covered in more detail in Chapter 6.²

```
1 def integer_identity(x: int):
2     if type(x) != int:
3         raise TypeError("x must be an integer")
4     return x
```

You can access function definitions and documentation using `?function_name` or `help(function_name)`.

¹See [common gotchas](#) from Reitz and Schlusser 2016 and Lubanovic 2019 page 147.

²See also the official documentation for raising exceptions.

5.5 Local and Global Variables and Namespaces

As we keep creating keyword parameters in our functions, or additional variables inside the function, you might start to wonder what if I had already used that keyword as a variable name? The below shows that no conflict arises. `foo` will not overwrite the variables `a` and `b`, and `foo` still operates without interference from those variables being previously defined.

```
1 a = 10
2 b = 10
3 def foo(a = 1):
4     b = 2*a
5     return b
6
7 x = foo(a = 0)
8 print(x)
9 print(a, b)
```

The reason this all behaves so well is that there is a global *namespace* in which `a` and `b` are defined as 10. And there is a separate namespace for the function `foo`, where `a` and `b` can exist independently, just as “die” can mean one thing in England and another thing in Germany.



Above: an argument over namespaces from *The Simpsons*.

There are local variables and there are global variables. A variable’s *scope* is the part of a program where the variable can be accessed. Variables created inside a function are local and their scope is the function. Different functions can have local variables of the same name without creating any kind of interference thanks to the separate namespaces.

5.6 Anonymous Functions

Anonymous functions are created as lambda functions. The main use for this is for functions that have a very short definition and that might only be used once in a program, or perhaps as an argument in a higher-order function (covered in Section 5.7).

```
1 def simple_function(x):
2     return 1 * (x > 0)
```

The above can be replaced by

```
1 simple_function = lambda x: 1 * (x > 0)
```

We can even do something a little more complicated.

```
1 simple_but_complicated = lambda x, y=-9: 1 * (x > y)
```

5.7 Higher-Order Functions Like map

A function that takes another function as an argument is a higher-order function (see Ramalho 2015). Two examples are the built-in functions `sorted` and `map`.

`sorted` takes an iterable and returns a new sorted iterable as long as every element of the iterable can be compared. Observe `sorted([3,1,2]) == [1, 2, 3]`. What makes `sorted` a higher-order function is the optional `key` argument.

`map` takes a function and applies it to every element in the iterable. It returns something call an iterator, holding the results. You can convert the iterator to a list using `list()` to see the results.

```
1 # Square the elements of a list
2 # [0,1,4,9]**2 doesn't work
3
4 f = lambda x: x**2
5 print(map(f, [0,1,2,3]))
6 print(list(map(f, [0,1,2,3])))
```

5.8 Examples

```
1 def letter_value(letter):
2
3     # convert to lowercase
4     letter = letter.lower()
5     # according to some blog
6     if letter in 'eaionrtlsu':
7         return 1
8     elif letter in 'dg':
9         return 2
10    elif letter in 'bcmp':
11        return 3
12    elif letter in 'fhvwy':
13        return 4
14    elif letter == 'k':
15        return 5
16    elif letter in 'jx':
17        return 8
18    elif letter in 'qz':
19        return 10
20    else:
21        return "not a valid letter"
22
23 def word_value(word):
24     return sum(map(letter_value, word))
```

Now run `sorted(['friends', 'enemies', 'burgers', 'I', 'a', 'I', 'I'], key = word_value)`.

5.8.1 Exercise!

```
1 # Find a nearby multiple of five
2 def find_close_multiple_of_five(num):
3     # check if multiple of five
4     is_multiple = num % 5 == 0
5     while is_multiple not True:
6         num += 1
7         is_multiple = num % 5
8     return num
```


Chapter 6

Exceptions

The [Zen of Python](#) advises us that “Errors should never pass silently.”

6.1 Exceptions

Reference: Lubanovic 2019 Chapter 9

Occasionally, you might ask Python to do something impossible. Try running `1/0`. You will get an error message, `ZeroDivisionError`. Robust code should deal with this possibility intelligently. This is especially true when writing functions. You can avoid errors by writing your code to prevent them from occurring. Or, you might write code that responds to errors. For more on the built-in exceptions, follow [this link](#).

Here’s an example of avoiding the division by zero error.

```
1 def pct_change(old, new):
2     delta = new - old
3     if old != 0:
4         pct = 100 * delta / old
5         return pct
6     else:
7         return "Not defined."
```

Here’s an example of dealing with the error, which requires `try` and `except` statements.

```
1 def pct_change1(old, new):
2     delta = new - old
3     try:
4         pct = 100 * delta / old
5         return pct
6     except ZeroDivisionError:
7         return "Not defined."
```

```
1 # Will this work?
2 def pct_change2(old, new):
3     delta = new - old
4     try:
5         pct = 100 * delta / old
6         return pct
7     return "Not defined."
```

A `try` must be paired with an `except`, so the definition of `pct_change2` will not work. There must be an `except` and, as we use in `pct_change1`, it’s a good idea to use the form `except ExceptionName`. This tells Python what code to run when a certain exception is encountered.

```
1 def pct_change3(old, new):
2     try:
3         delta = new - old
4     except TypeError:
```

```

5         return "Use ints or floats."
6     try:
7         pct = 100 * delta / old
8         return pct
9     except ZeroDivisionError:
10        return "Not defined."

```

The following function, `pct_change4`, actually won't run properly if you attempt to execute `pct_change4(1, 'cheese')`. Can you figure out why? Think about the ordering of the code.

```

1 def pct_change4(old, new):
2     delta = new - old
3     try:
4         pct = 100 * delta / old
5         return pct
6     except ZeroDivisionError:
7         return "Not defined."
8     except TypeError:
9         return "Use ints or floats."

```

You do not need to specify the error type in your `except` statement. Sometimes this might hide errors that you do want to be surfaced or stop the execution of a program, so use these blanket exceptions carefully.

```

1 def pct_change5(old, new):
2     try:
3         delta = new - old
4         pct = 100 * delta / old
5     except:
6         return "An error occurred."
7     return pct

```

You might use a blanket `except` after handling specific errors.

```

1 def pct_change6(old, new):
2     try:
3         delta = new - old
4         pct = 100 * delta / old
5         return pct
6     except ZeroDivisionError:
7         return "You can't divide by zero."
8     except:
9         return "An error occurred."

```

It can be useful information to know what type of exception your code generates. In that case, you can access and print that error.

```

1 try:
2     'a' + 1
3 except Exception as e:
4     print(e)

```

The use of `Exception` above matters. Compare these two programs.

```

1 try:
2     'a' + 1
3 except TypeError as e:
4     print(e)

```

```

1 try:
2     'a' + 1
3 except ValueError as e:
4     print(e)

```

Only the first program above actually handles the exception.

Now, we consider the use of `else` and `finally` statements after a `try/except`. An `else` clause can be added after the `except` clauses, and the statements in the `else` clause are executed only if no exceptions were raised. The `else` is then in contrast to the raising of exceptions, in a similar style as when an `else` might be used in complement to an `if`. What output do you expect from the following program? What if we changed the value of `denom`?

```
1 denom = 0
2 try:
3     print(2 / denom)
4 except Exception as e:
5     print(e)
6 else:
7     print(3 / denom)
```

Here are slightly more complicated examples. Try running them to see what happens.

```
1 products = ['bike', 'treadmill']
2 try:
3     print(products[2])
4 except IndexError as e:
5     print(e)
6 except Exception as e:
7     print(e, 'other error')
8 else:
9     print(products[2], '!')
```

```
1 products = ['bike', 'treadmill']
2 try:
3     print(products[1])
4 except IndexError as e:
5     print(e)
6 except Exception as e:
7     print(e, 'other error')
8 else:
9     print("Else block time")
10    print(products[2], '!')
```

Next, we introduce the `finally` clause. Whereas the `else` block was executed when no exceptions were raised, the `finally` block is executed no matter what. This [StackExchange post](#) helps explain the unique usefulness of this, but we won't get into that much detail.

```
1 products = ['bike', 'treadmill']
2 idx = 2
3 try:
4     print(products[idx])
5 except Exception as e:
6     print(e, 'other error')
7 finally:
8     print('End')
```


Chapter 7

More Data Structures

7.1 Tuples

Reference: Lubanovic 2019 Chapter 7

A *tuple* is a lot like a list. Whereas lists used square brackets (`[]`), a tuple uses parentheses, (`()`). Like lists, we can index and slice a tuple. The main difference is that tuples are immutable. We cannot reassign the element at a particular index.

```
1 example_list = ['bike', 'treadmill']
2 example_tuple = ('bike', 'treadmill')
3
4 example_list[1] = 'treadmill'
5 example_tuple[1] = 'spacecraft' # This will throw an error
```

Note we can convert lists and tuples, similar to the way we could convert floats and strings (`int("1")`, `str(1)`).

```
1 example_list = ['bike', 'treadmill']
2 example_tuple = ('bike', 'treadmill')
3
4 print(list(example_tuple))
5 print(tuple(example_list))
```

It might be hard to see immutability as an advantage, but this does make tuples to be safer objects. You won't accidentally screw them up (unless at the very beginning). There is another more definite advantage to tuples. They are processed faster. Processing speed will not be a practical concern in this class, but it could be a concern in your professional career.

7.1.1 Tuple Assignment

Tuples can also be created without parentheses.

```
1 friends = 'Big Bird', 'Snuffleupagus'
```

Or you can we can create individual variable names with something like the following.

```
1 main, sidekick = 'Mitt', 'Paul'
```

This is kind of nice and or tricky when thinking about reassigning variables. Compare the following.

```
1 apples = 10
2 pies_possible = 2
3
4 apples, pies_possible = 2*apples, apples/5
5 print(apples, pies_possible)
6
7 apples = 10
8 pies_possible = 2
9
```

```

10 apples *= 2
11 pies_possible = apples/5
12 print(apples, pies_possible)

```

7.1.2 Unpacking

Tuples can be *unpacked* with an asterisk, `*`, to create a starred expression. This is useful when the tuple structure isn't wanted, like when you don't want nested tuples.

For example, let's start with a tuple, `friends = 'Antonin', 'Ruth'`. We might want to add these to a larger tuple. Let's try:

```
1 former_justices = 'Breyer', friends
```

This results in a nested tuple `('Breyer', ('Antonin', 'Ruth'))`. We should unpack `friends`.

```
1 former_justices = 'Breyer', *friends
```

Now we have a flat tuple `('Breyer', 'Antonin', 'Ruth')`. This would also work if `friends` were a list. That is, we get the same value for `former_justices` if we instead ran `former_justices = 'Breyer', *list(friends)`.

Next, unpacking can be useful in use with functions when you'd like each element of the tuple to be processed as a new positional argument.

```

1 def heart(a,b):
2     return a + " <3 " + b
3 # heart(friends) returns an error
4 print( heart(*friends) )

```

7.1.3 Tuples vs. Lists

Recall the primary difference between lists and tuples: lists are mutable and tuples are not. It might be hard to see immutability as an advantage, but this does make tuples to be safer objects. You won't accidentally screw them up (unless at the very beginning).

There is one definite advantage to tuples. They are processed faster. Processing speed will not be a practical concern in this class, but it could be a concern in your professional career.

7.2 The `in` Operator

The `in` operator allows you to determine if an element is contained in a list or tuple. In Python, the statement `x in A`, where `A` is a list, mirrors the mathematical statement $x \in A$ where A is a set. While we never formally introduced `in`, we've already seen it in for loops with the for clause, `for item in range(10):`, for example.

Realize that for numerics, `in` will evaluate to true as long as the element is equal (`==`) to something in the list. That means a float can be `in` a list of integers.

```

1 chapters_on_midterm = [2,3,4,4,5,6,7,8,9]
2
3 # confirm all integers
4 for item in chapters_on_midterm:
5     print(item, type(item))
6
7 # This is certainly True
8 bool1 = 2 in chapters_on_midterm
9
10 # What about this?
11 bool2 = 2.0 in chapters_on_midterm # This is True!

```

Finally, you can use `not in` as you might expect. The statement `x not in some_list` if there is no element `y in some_list` that is equal to `x`.

7.3 List Methods

Reference: Lubanovic 2019 Chapter 7

Python methods are like functions, but they are associated with objects.

Functions look like this: `sorted(some_list)`

Methods look like this: `some_list.sort()`

For now, we can proceed thinking of them just as functions with this syntax.

The function `sorted()` and the method `.sort()` do the same thing in some sense—both can be used to sort a list. They are different in that `sorted()` returns a new sorted list without mutating `some_list`. The method `.sort()` doesn't return anything; it mutates the list into a sorted list. This difference isn't a property of functions and methods. The difference is particular to `sorted()` and `.sort()`.

```
1 # Demonstration of sorted() and .sort()
2
3 presorted_list = [1,2,3,4]
4 alt_list = [1,4,2,3]
5
6 c1 = alt_list == presorted_list
7 print(c1)
8
9 c2 = sorted(alt_list) == presorted_list
10 print(c2)
11
12 c3 = alt_list == presorted_list
13 print(c3)
14
15 alt_list.sort()
16
17 c4 = alt_list == presorted_list
18 print(c4)
```

Other important list methods include `.append()` and `.index()`. `.append()` requires an argument—an element to be added to the end of the list.

Running `some_list.append(x)` does the same thing `some_list += [x]` would accomplish.

```
1 # .append() Demonstration
2 ones = list()
3
4 ones.append(1)
5
6 print(ones)
7
8 ones += [1.0] # Recall this is the same as ones = ones + [1.0]
9
10 print(ones)
```

Next, the `.index()` method helps us find the index of the first instance of a particular element in a list. This method requires an argument, and `some_list.index(x)` returns the minimum index of `x` in `some_list`.

```
1 # .index() Demonstration
2 ones = [1, 1.0]
3
4 print(ones.index(1))
5
6 print(ones.index(1.0))
7
8 # There is no distinction between int and float
9 print(ones.index(1.0) == ones.index(1))
```

Exercise: Remove the duplicates from a list.

```
1 big_list = [1,2,4,2,12,4,12,234,1,1,1,1] # a lot duplicates
2
3 big_list_without_dups = []
4
```

```

5 for element in big_list:
6
7     if element not in big_list_without_dups:
8
9         big_list_without_dups.append(element)
10
11
12 ## Alternate version
13
14 big_list_without_dups_alt = []
15
16 idx = 0 # keep track of index of each element
17 for element in big_list:
18     # Add the element if it's the first time we've seen it in big_list
19     first_index = big_list.index(element)
20     if idx == first_index:
21         big_list_without_dups_alt.append(element)
22
23     # Advance the index for the next element
24     idx += 1

```

7.4 Strings

Reference: Lubanovic 2019 Chapter 5

Strings behave like lists in terms of indexing, slicing, and iterating. Strings are immutable however.

The `in` operator can be used on strings to test if one string is a substring within another string.

```

1 for char in 'team':
2     print(char, char in team):
3
4 print("I" in 'team')

```

Important string methods include `.upper()`, `.lower()`, `.isalpha()`, `.split()`, and `.replace()`.

The `.upper()` and `.lower()` methods return a new string in uppercase or lowercase letters, respectively. These do not mutate the original string and no arguments are necessary.

The `.isalpha()` method returns a boolean, stating whether or not the string contains all alphabetical characters (this is different than not being an integer data type for example—a string might still contain a numeric character).

```

1 # In case you can't trust your eyes for l vs 1.
2
3 lowercase_L = "l"
4 one = "1"
5
6 for item in lowercase_L, one:
7     print(item, item.isalpha())

```

The `.replace()` method takes two arguments. It returns a new string that replaces every instance of the first argument with the second argument. The function below returns a more muted string by eliminating all capital letters and converting exclamation marks to periods.

```

1 def lower_your_voice(string):
2     lowercase = string.lower()
3     not_exclamatory = lowercase.replace("!", ".")
4     return not_exclamatory
5
6 print(lower_your_voice("Your card was declined!"))

```

The `.split()` method requires an argument and returns a list, dividing a string into substrings based on the argument. This is especially useful for splitting a sentence into its individual words.

```

1 invisible_hand = "It is not from the benevolence of the butcher that we expect our dinner
2     but from their regard to their own interest"
3
4 the_words = invisible_hand.split(" ")

```



```

4
5 print(the_words)
6
7 # Note what happens when you split on the first or last character in a string
8
9 laugh = 'hahahahah'
10
11 split_laugh = laugh.split('h')
12
13 print(split_laugh)
14
15 # Note what happens when you pass no argument
16 print(laugh.split())
17
18 # This gives an error. You can't split on a zero-length separator.
19 laugh.split("")

```

7.5 Dictionaries

Reference: Lubanovic 2019 Chapter 8

Dictionaries allow us to store key-value pairs. It's kind of like having one row in a table. Keys are like column names and values are the row-column cell value. Key-value pairs are created as `key: value`, then separated by commas and wrapped in curly braces, `{}`, to create a dictionary. Consider the example below.

```
1 workout = {'user': 'Velma', 'fitness_discipline': 'cycling', 'instructor': 'Matt Wilpers'}
```

A specific value is access by indexing the dictionary by the key.

```
1 print(workout['user'])
```

We can add new key-value pairs by assigning the value to the dictionary at that key.

```

1 # Build a dictionary from scratch
2 journal = dict() # creates an empty dictionary, can also use {}
3
4 journal['2020-10-03'] = "Today I learned a lot of Python. It was buckets of fun."

```

The `in` operator works on dictionaries by searching the keys.

```

1 # Help translate bad journalism
2 media_translator = {'is caused by': 'is correlated with'}
3
4 print('is caused by' in media_translator)
5 print('is correlated with' in media_translator)

```

You can access the keys with the `.keys()` operator and values with the `.values()` operator. So `x in some_dict` is actually a shorthand for `x in some_dict.keys()`.

Dictionary keys must be immutable. Tuples are fine. Lists are not.

```

1 # Economist Santa
2
3 gifts = {} # could also use dict() here
4
5 for child in ['Anna', 'Boris']:
6     for year in [2020, 2021]:
7
8         key = child, year # this is a tuple just like (child, year)
9
10        gifts[key] = 'money'
11
12 print(gifts)

```

7.6 Sets

Reference: Lubanovic 2019 Chapter 8

I find sets to be underrated. They give flexibility in analysis because they allow for quick intersections and unions (e.g. find and analyze users who did this *and/or* that).

Sets are unordered and duplicates are ignored. We construct them like lists, but use `{}` instead of `[]`. Note that to create an empty set, you should use `set()` because `{}` creates an empty dictionary.

Being unordered means it is true that `{1,2} == {2,1}`.

It's not exactly right to say that sets *can't* have duplicates. You can create a set with duplicate elements and no error will be thrown. But those duplicates are ignored so that the created set object will not in fact have any duplicates. Thus, it is true that `{1,2} == {2,1,1,2}`.

Three important methods are `.union()`, `.intersection()`, and `.difference()`. Each of these acts on a set and requires another set as an argument. Intersection and union work like the set operations \cap and \cup . The difference method performs set subtraction, \setminus . Recall that set subtraction is not commutative; $A \setminus B \neq B \setminus A$ unless $A = B$.

```

1 # Union and Intersection
2
3 primes = {2,3,5}
4 evens = {2,4,6}
5
6 even_and_prime = primes.intersection(evens)
7
8 even_or_prime = primes.union(evens)
9
10 for set_ in even_and_prime, even_or_prime: # note we're iterating over a tuple
11     print(set_)
12
13 # Set Subtraction
14
15 contiguous_USA = {'New York', 'Kentucky', 'Wisconsin', 'California'} # among other states
16
17 tectonic_seceder = {'California'}
18
19
20 print(contiguous_USA.difference(tectonic_seceder))
21
22 # We can also use the - operator
23 print(contiguous_USA - tectonic_seceder)
24
25 ## More Subtraction
26
27 cold_places = {'Wisconsin', 'Yukon'}
28
29 print(contiguous_USA.difference(cold_places))
30
31 # Reverse the arguments
32 print(cold_places.difference(contiguous_USA))

```

Now that we've seen these data types, you are well prepared to work with a lot of data you might find in the wild. See the code from class for an example using a public API.

Chapter 8

Modules

Reference: Lubanovic 2019 Chapter 11

A *module* is a file that contains Python code. Large programs are more manageable when divided into modules. Many functions in the standard Python library are stored in modules. The `math` and `random` modules are common examples. These shouldn't require additional installation to use if you've downloaded Anaconda.

To use a module, you must import the module with an import statement, `import math` for example. Then any function in that module can be accessed by using the function name, prefixed by the module name and a dot. To use `sqrt` from `math`, you must use `math.sqrt(81)`. [PEP 8](#) advises, "imports are always put at the top of the file, just after any module comments and docstrings, and before module globals and constants."

You can import just a specific function with syntax like the following: `from math import sqrt`. Then `sqrt` can be accessed without the `math.` prefix. Running `from math import *` is called a wildcard import and will import every function in the module.

Finally, you can *alias* a specific module, library, or function with an `as` clause. There are conventional aliases for many modules. Pandas, NumPy, and Datetime are libraries we will cover later. These are typically imported with aliasing as in the below.

```
1 import pandas as pd
2 import numpy as np
3 import datetime as dt
4 from math import sin as sine # Not a typical alias, but for demonstration
```

Note that if you import a module inside a function, it's not available globally. This is similar to how local and global variables work.

```
1 def silly_function():
2     import numpy as np
3
4 silly_function()
5
6 # Next line produces an error. NumPY is not imported.
7 np.array([])
```

8.1 Storing Functions in Modules

You can create your own module by placing code into a `.py` file. If that file is in your directory, you can access it with an import statement. In the lecture folder, I've placed a file `next_power.py`. Download that and place it in your working directory to try importing it. In Google Colab, you can upload files in the left menu.

Try the following.

```
1 import next_power
2 val = next_power.next_power_of_five(126)
3 print(val)
```

```
1 import next_power as npow
2 val = npow.next_power_of_five(44)
3 print(val)
```

```
1 from next_power import next_power_of_five
2 val = npow.next_power_of_five(309)
3 print(val)
```

It's difficult to overstate how helpful this is in creating cleaner and more readable Jupyter notebooks.

8.2 Using the random module and matplotlib

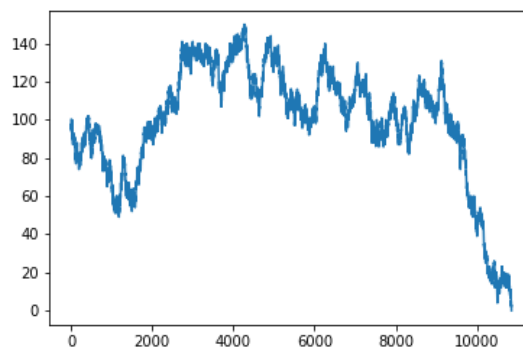
Let's get our hands dirty a bit.

```
1 import random
2 random.seed(33) # pick a seed for reproducibility
3
4 number = random.randint(0,1)

1 # Gambler's Ruin
2
3 purse = 100 # starting money
4
5 # Keep gambling if you have money
6 total_gambles = 0
7 purse_values = [purse]
8 while purse > 0:
9     outcome = random.choice([-1,1]) # win or lose 1 with eve odds
10    purse += outcome
11
12    purse_values.append(purse)
13    total_gambles += 1
```

We can graph this with matplotlib.

```
1 import matplotlib.pyplot as plt
2 plt.plot(range(total_gambles+1), purse_values)
3 plt.show()
```



Chapter 9

Object Oriented Programming

9.1 Object Oriented Programming

Reference: [Lubanovic 2019 Chapter 10](#)

So far, our programming has been *procedural*. A program was made of procedures and data might be passed from one procedure to the next. This creates a separation of data and the procedures/operations. As a program grows, this can become more unwieldy. We'll talk about defining classes. While you might get quite far without having to define your own classes, the vocabulary here is important for anyone who wants to be fluent in Python and understand popular libraries like Pandas.

Object-oriented programming (OOP) centers on creating objects instead of procedures. An object contains both data and procedures. An object's data are its *data attributes*. An object's procedures are its *methods*, which operate on the data attributes. Attributes are like variables and methods are like functions. Bundling data with the code operating on it is called *encapsulation*.

9.1.1 Classes

A *class* is code that specifies the data attributes and methods for a particular type of object. A class is like a blueprint and the object is the particular realization. When we previously talked about data types, we were really referencing classes. Run `print(type('Hello, World!'))`. Your output should be `<class 'str'>`.

Classes are defined in the following way. Note, per [PEP 8](#), class names should follow the CapWords convention.

```
1 import random
2
3 class Coin:
4
5     def __init__(self):
6         self.sideup = "Heads"
7         self.coin = "Quarter"
```

Run `coin = Coin()` and print `coin.sideup` and `coin.coin`.

So a minimal class definition includes the `__init__` *initializer method*. This initializes the objects data attributes. We could instead make these attributes arguments.

```
1 import random
2
3 class Coin:
4
5     def __init__(self, sideup = 'Heads', coin = 'Quarter'):
6         self.sideup = sideup
7         self.coin = coin
```

Take care to note the last two lines in the block above are necessary to create the `sideup` and `coin` attributes. Now, let's add some methods to our class.

```

1 import random
2
3 # Simulate a coin that can be tossed
4
5 class Coin:
6
7     def __init__(self, sideup = 'Heads', coin = 'Quarter'):
8         self.sideup = sideup
9         self.coin = coin
10
11     def toss(self):
12         toss_outcome = random.choice(['Heads','Tails'])
13         self.sideup = toss_outcome # Here we change the attribute instead of returning a
    value

```

See what happens now.

```

1 coin = Coin() # start with a coin that is heads up per class default value
2 for i in range(100):
3
4     print(coin.sideup)
5     coin.toss() # toss the coin

```

Here we add a value-returning method.

```

1 import random
2
3 # Simulate a coin that can be tossed
4
5 class Coin:
6
7     def __init__(self, sideup = 'Heads', coin = 'Quarter'):
8         self.sideup = sideup
9         self.coin = coin
10
11     def toss(self):
12         toss_outcome = random.choice(['Heads','Tails']) # local variable just like before in
    defining functions
13         self.sideup = toss_outcome # Here we change the attribute instead of returning a
    value
14
15     def get_sideup(self):
16         return self.sideup

```

Next, we might want to use *hidden attributes*. Before we could externally change `sideup` attribute. Perhaps you don't want that to be possible in this or in another setting. Then you can make that attribute private by including two underscores with the init.

```

1 import random
2
3 # Simulate a coin that can be tossed
4
5 class Coin:
6
7     def __init__(self, sideup = 'Heads', coin = 'Quarter'):
8         self.__sideup = sideup
9         self.coin = coin
10
11     def toss(self):
12         toss_outcome = random.choice(['Heads','Tails']) # local variable just like before in
    defining functions
13         self.__sideup = toss_outcome # Here we change the attribute instead of returning a
    value
14
15     def get_sideup(self):
16         return self.__sideup # We need two underscores here too!

```

Now for `coin = Coin()`, you cannot access the `sideup` attribute with either `coin.sideup` or `coin.__sideup`. The `get_sideup()` method becomes necessary to access the private attribute. When making a data attribute

private, one might create methods for accessing and changing those attributes. These are called accessor and mutator methods. Or you might call them getters and setters, respectively.

The `__str__` method is designed to indicate an object's emphstate (the attribute values).

```

1 import random
2
3 class Coin:
4
5     def __init__(self, sideup = 'Heads', coin = 'Quarter'):
6         self.sideup = sideup
7         self.coin = coin
8
9     def __str__(self):
10        return "The coin is a " + self.coin + ", and it is " + self.sideup + "."

```

This method is accessed not directly, but by printing the object.

9.2 Inheritance

Inheritance allows a new class to extend an existing class. This helps with code reusability a bit because we can have super and subclasses.

We might start with a superclass. Here's an example.

```

1 class Automobile():
2
3     def __init__(self, gas_tank):
4         self.gas_tank = gas_tank
5
6     def drive(self):
7         self.gas_tank -= 1

```

Now let's make a subclass.

```

1 class HybridCar(Automobile):
2
3     def __init__(self, gas_tank, battery):
4         Automobile.__init__(self, gas_tank)
5
6         # initialize additional battery parameter
7         self.battery = 100

```

Try defining `prius = HybridCar(50,100)` and running `prius.drive()`. It should work. Check `prius.gas_tank`. This illustrates basic inheritance of classes.

9.2.1 Polymorphism (G§11.2)

Now we demonstrate *polymorphism*. A subclass can have the same methods defined as their superclass. The methods override the superclass.

```

1 class HybridCar(Automobile):
2
3     def __init__(self, gas_tank, battery):
4         Automobile.__init__(self, gas_tank)
5
6         # initialize additional battery parameter
7         self.battery = battery
8
9     def drive(self):
10        self.gas_tank -= .5
11        self.battery -= 1

```

Python includes a handy `isinstance()` function that helps determine if an object is of a certain class or of an instance of a subclass of that class.

```

1 print(isinstance(prius, Automobile))
2 print(isinstance(prius, HybridCar))

```


Chapter 10

IPython

10.1 IPython and Jupyter

Reference: McKinney 2017 Chapter 2 and VanderPlas 2016a Chapter 1

Thus far, we've been using Google Colab, which runs IPython notebooks. Now, we'll use Jupyter. Basically, we're moving off of the cloud.

10.1.1 IPython Basics

The Jupyter notebook is an interactive document for code, text, data visualization, and other output. Python's Jupyter kernel uses the IPython system, so all of the IPython basics we'll cover will apply to Jupyter notebooks. Beyond simply executing Python code, IPython comes with enhanced features that make coding easier. I've highlighted a few here. Read through McKinney §2.2 or check the IPython sections from Vanderplas's [Python Data Science Handbook](#).

Tab completion is a feature that allows you to hit tab following some input and then any variables matching the characters will be displayed.

Introspection refers to the ability to place a `?` before or after a variable and executing the line will display helpful information about the object. For a function, you'll be shown the *docstring*.

Notebooks are divided into discrete code cells. You can select multiple with shift and select. Then use **shift-m** to **merge cells**. Use **ctrl-c** to interrupt a command.

There are also certain **magic commands** that begin with a `%`. My favorite are `%timeit` and `%load`. See McKinney page 29 for a table of magic commands or Vanderplas [here](#). Another useful magic command when working with your own modules is the [autoreload command](#). If you are developing in a separate `.py` file and prototyping in a notebook, this is very handy.

```
1 %load_ext autoreload
2 %autoreload 2
```

You can run **shell commands** with `!` or `%` as well. The former creates a subshell, so it cannot be used for changing your working directory. See more [here](#) (Vanderplas).

10.1.2 Extensions

I have installed the *code-folding* extension, which I recommend to anyone who might be dealing with messy notebooks.

10.1.3 Peak Under the Hood: Grading

This is on the more complicated end, but here is how I intended to grade your `.py` files from the midterm.

```
1 # all submissions are store in folder PyFiles
2 import os # module for additional string-based shell commands
3 from importlib import reload # allows for reloading a module
4
5 # ls lists all files in a folder
6 # Put all file names in a list
7 files = !ls ~/PyFiles
8
9 # Change working directory
10 %cd ~/PyFiles
11
12 for file in files:
13     # rename the file
14     os.system("mv " + file + " temp.py")
15     import temp
16     reload(temp) # overwrites previous iteration import
17
18     # grade file
19     print(file)
20     print(temp.add(2,3)) # all files contain a function named add
21
22     # name the file back
23     os.system("mv temp.py " + file)
```

Part II

Data

Chapter 11

NumPy

Reference: McKinney 2017 Chapter 4, VanderPlas 2016a Chapter 2

NumPy (Numerical Python) is an important library for working with arrays. This shouldn't require any additional installation in Anaconda and is available in Google Colab. NumPy is conventionally imported with the alias `np`: `import numpy as np`.

11.1 The Array

The NumPy array is important.

```
1 import numpy as np
2
3 empty_array = np.array([])
4 print(len(empty_array))
5
6 zero1 = np.array([0])
7 print(len(zero1))
8
9 zero2 = np.array(0)
10 #print(len(zero2))
11
12 #try_this = np.array(1,2)
13
14 two_d = np.array(((1,2),(3,4)))
15
16 for array in [empty_array, zero1, zero2, two_d]:
17     print(type(array))
```

The array is of a class `ndarray`.

Finally, note arrays can be of mixed type. Mixed type can be constructed by specifying the data type as `object`, `np.array(['s',1], dtype = object)`. Still I don't know a good reason to use NumPy with mixed types. And as the name suggests, NumPy is designed for numeric data. Notice that `np.array(['s',1])`, without specifying `dtype = object`, results in the `1` being converted to a string.

11.2 Indexing

Basic indexing is done like for lists. However, arrays can have multiple dimensions.

```
1 arr = np.arange(10)
2
3 print(arr[5]) # the 5th index
4
5 # slice
6 print(arr[5:8])
7 # reassign
8 arr[5:8] = 0
```

```

9 print(arr)
10
11 print(arr[10]) # error remember we have zero-based indexing

```

Array slices are *views*. The data is not copied and any modifications are transferred to the source array.

```

1 arr_slice = arr[5:8]
2 print(arr_slice)
3
4 arr_slice = 1 # does nothing bc it writes over the slice with an int
5 print(arr, arr_slice)
6
7 arr_slice = arr[5:8]
8 print(arr, arr_slice)
9
10 arr_slice[:] = 1, 2, 3 # mutates both
11 print(arr, arr_slice)
12
13 arr_slice[:] = -40 # mutates both
14 print(arr, arr_slice)

```

With multidimensional arrays, there's a bit more to indexing.

```

1 arr2d = np.array( ( (1,2), (8,9) ) )

```

Thinking of the two-D array as a matrix, the first index place will select the row and the second will select the column.

```

1 print(arr2d[0])
2
3 print(arr2d[0][0])
4 print(arr2d[0,1]) # these are the same
5
6
7 print(arr2d[:,-1]) # get the last item from each row

```

Check Figure 4-2 in McKinney for a good illustration of slicing two-D arrays.

Boolean Indexing

Recall `arr` from above. We left it with value `arr = np.array([0, 1, 2, 3, 4, -40, -40, -40, 8, 9])`. We can index it based on a boolean condition using a syntax `arr[<condition>]`

```

1 print(arr)
2
3 print(arr[arr > 0]) # reduces the array

```

We can combine conditions with logical negations, ands and ors, but *not* with the `not`, `and`, `or` keywords. Use `~` to negate an array of booleans, `&` for and, and `|` for or.

11.2.1 Functions

Universal functions perform element-wise operations.

```

1 arr = np.arange(10)
2
3 print(np.sqrt(arr))

```

There are also import statistical functions like `np.mean()` and `np.std()` for averaging and finding the standard deviation.

11.2.2 Linear Algebra

We can access linear algebra functions using the `numpy.linalg` submodule.

As demonstrated above, operators like `*` apply element wise. If we have an array `x = np.array([1,2])`, then `x * x` gives an output of `np.array([1,4])`. Thus, `*` is not the dot product \cdot . To get the scalar-valued $x \cdot x = \sum x_i^2$, we need `np.dot(x,x)`. The below illustrates the same with matrix multiplication.

```

1 x = np.array([[1,0],[0,1]]) # identity
2 y = np.array([[8,1],[2,3]])
3
4 mystery1 = x * y
5 mystery2 = np.matmul(x,y)

```

You'll find `mystery2 == y` which is what should be expected for matrix multiplication. `mystery1` gives array `[[8, 0],[0, 3]]`, meaning element-wise multiplication was done. Python didn't know you wanted matrix multiplication. You can however call create a specific matrix object with `np.matrix()`. Let's repeat the above.

```

1 x = np.matrix([[1,0],[0,1]]) # identity
2 y = np.matrix([[8,1],[2,3]])
3
4 mystery1 = x * y
5 mystery2 = np.matmul(x,y)
6
7 print(mystery1 == mystery2)

```

You can invert a matrix with `np.linalg.inv()`, and it will accept either an array or `np.matrix` object.

Regression Exercise

Let's illustrate with a simple linear regression with no intercept. Suppose we have a single independent variable and linear model

$$y = \beta x + \epsilon.$$

Let's simulate some data. We'll assume $\beta = 1$, $\epsilon \sim N(0,1)$, and find $\hat{\beta}$ from simulated data. Recall $\hat{\beta} = (X^T X)^{-1} X^T Y$ so this is a matter of matrix multiplication.

```

1 import numpy as np
2
3 n_obs = 100 # observations
4 x = np.random.random(n_obs)
5 epsilon = np.random.normal(0,1,n_obs)
6 true_beta = 1
7
8 y = true_beta * x + epsilon
9
10 # make 'tall' matrices with rows for each observation
11 x_matrix = np.matrix(x).T
12 y_matrix = np.matrix(y).T
13
14 # Check on linear algebra and operators
15 (x_matrix.T * x_matrix)**-1 == np.linalg.inv( np.matmul(x_matrix.T, x_matrix) )
16 # above actually compares element to element
17
18
19 beta_hat = (x_matrix.T * x_matrix)**-1 * (x_matrix.T * y_matrix)

```

You can find more linear algebra examples on my [github](#).

Chapter 12

Pandas: Series and DataFrames

Pandas was created by Wes McKinney. If you'll work with data in Python, get used to having this at the top of your code.

```
1 import pandas as pd
```

12.1 Series

A Pandas *series* is like an array. Unlike lists or numpy arrays, a series has an *index*. If you create a series from scratch, the default index will be numbered from zero. You can also explicitly pass an index. These are constructed much like a numpy array, but with `pd.Series()`. You can also pass a dictionary to create a series, where the keys are the index values. You can go back to a NumPy array with the method or `.to_numpy()` by accessing the values attribute, `.values`.

```
1 ser1 = pd.Series([10,10,2,20])
2 ser2 = pd.Series(['squirrel', 100], index = ['animal', 'number'])
3
4 dictionary_helper = {'Today': 60, 'Tomorrow': 50, 'Monday': 55}
5 ser3 = pd.Series(dictionary_helper)
6
7
8 print(ser1)
9 print(ser2)
10 print(ser3)
11
12 # Compare to Numpy
13 import numpy as np
14 print(np.array(ser1)) # no more index
15 print(np.array(ser2))
16
17 # Compare ser3 to its dictionary
18 print("Today" in dictionary_helper)
19 print("Today" in ser3)
20
21 print(60 in dictionary_helper)
22 print(60 in ser3)
```

We can access the series index and values similarly as we'd access a dictionary's keys and values.

```
1 print(ser3.index)
2 print(ser3.values)
```

Mathematical operations are automatically aligned based on the series index.

```
1 wealth = {'Alice': 65, 'Bob': 50}
2 wealth_series = pd.Series(wealth)
3
4 bonus = {"Bob": 10, "Alice": 10}
```

```

5 bonus_series = pd.Series(bonus)
6
7 print(wealth_series + bonus_series)
8
9 # convert to numpy and add
10 print( np.array(wealth_series) + np.array(bonus_series) )

```

Null values can arise in a series. An index might be explicitly specified and without any corresponding value, or an operation might not be possible for a particular index.

```

1 wealth = {'Alice': 65, 'Bob': 50}
2 wealth_series = pd.Series(wealth, index = ['Larry', 'Alice', 'Debbie', 'Bob'])
3
4 bonus = {"Bob": 10, "Alice": 10}
5 bonus_series = pd.Series(bonus)
6
7 new_wealth = wealth_series + bonus_series
8 print(new_wealth)

```

There are a few methods to help with nulls. The methods `isnull()` and `notnull()` return booleans as the names would suggest.

```

1 print(new_wealth.isnull())
2 print(new_wealth.notnull())
3
4 print(new_wealth.isnull() | new_wealth.notnull()) # guess what this will be

```

You can access the

12.1.1 Series Functionality I

Like with NumPy arrays, you can add two series, multiply two series, etc. You can also add, multiply, etc by a constant.

```

1 a = bonus_series + 1
2 b = bonus_series + bonus_series
3 c = bonus_series * 2 * bonus_series

```

You can also apply NumPy functions that will operate element-wise.

```

1 a = np.sqrt(bonus_series)
2 b = np.exp(bonus_series)
3 c = np.log(b)

```

Apply

Some functions don't automatically apply to sequences element-wise, but you might want them to. The `apply()` method is made for these cases. Pass a function to `apply` and it will be applied at every index in the series.

```

1 def odd_or_even(x):
2     if x % 2 == 0:
3         return "Even"
4     return "Odd"
5
6 ser = pd.Series(range(1,9))
7
8 # This gives an error
9 odd_or_even(ser)
10
11 # Use apply
12 ser.apply(odd_or_even)

```

Anonymous functions are especially useful with the `apply` method.

Anonymous Functions

Anonymous, or lambda, functions are defined without the `def` keyword. They are nameless and can come in handy when needed for a short period. They are often used inside other functions. They follow a syntax like `lambda [argument]: [expression to return]`.

```
1 example = lambda x: x+1
2 print(example(-1))

1 ser.apply(lambda x: '2' in str(x))
```

Accessor Methods

Better than using `apply()` might be using accessor methods.

The string accessor method works by including `.str` after a Series and then a string method.

```
1 df.string_column.str.lower()
2
3 # better than
4 #df.string_column.apply(lambda x: x.lower())
```

12.2 DataFrames

You can imagine a series with multiple columns. That would be a dataframe, `pd.DataFrame`. Below are a few constructions.

```
1 # construct some DataFrames()
2 a = pd.DataFrame() # empty
3 b = pd.DataFrame(ser)
4 c = pd.DataFrame(ser, ser) # less common
5 d = pd.DataFrame([ser,ser]) # less common
```

DataFrames are also commonly constructed with a dictionary.

```
1 data = {'State' : ['KY', 'NY'],
2         'Capital' : ['Frankfort', 'Albany']}
3 df = pd.DataFrame(data)
```

You can also read a CSV with `pd.read_csv()`.

```
1 atus_df = pd.read_csv("ATUS_activity_2019.csv")
```

The `head()` and `tail()` methods to display the first or last rows. By default, five rows will be selected.

A specific column can be accessed with dict-like notation or by attribute. As shown in Antao 2022, dictionary access can be faster.

```
1 df['State']
2 df.State
```

Similarly, a new column can be created with the same dict-like notation.

```
1 df['Extra Column'] = None
```

Rows and columns can be dropped with the `drop()` method. Columns can also be deleted with `del`.

```
1 df['Extra Column'] = None
2 print(df.columns)
3
4 del df['Extra Column']
5 print(df.columns)
6
7 df['Extra Column'] = None
8 df.drop('Extra Column', axis = 'columns', inplace = True)
9
10 # axis = 1 also references columns
11 df['Extra Column'] = None
12 df.drop('Extra Column', axis = 1, inplace = True)
```

12.2.1 Indexing

You can specify an existing column as the index with `set_index()`. You can also explicitly change the index by accessing the index attribute. You can reset the index with `reset_index()`.

```
1 print(df.index)
2
3 df.index = [1, 'clown']
4
5 df.set_index("State") # returns a new dataframe
6 df.set_index("State", inplace = True) # alters df
7
8 # go back to numbered index
9 df.reset_index(inplace = True)
```

A DataFrame can be index with either `loc` or `iloc`. Use `loc` to index by the exact index and column names. Use `iloc` to index by the index and column numbers. An index can contain duplicates, which can complicate the below.

```
1 # return to State index
2 df.set_index('State', inplace = True)
3
4 a = df.loc['KY', 'Capital']
5 b = df.iloc[0, 0]
6
7 print(a,b)
```

As you could select an entire column with `df['Capital']`, you can select an entire row with `df.loc['KY']` or `atus_df.loc[0]`. Or you can select a subset by passing a list or slicing

```
1 sub_df1 = atus_df.loc[[0,10,29]]
2 sub_df2 = atus_df.loc[0:2]
```

12.3 Summarizing and Computing Descriptive Stats

Let's return to `atus_df`. Let's examine sleep averages and find the person who slept for the longest.

First we can mask to just select the rows where the activity is sleeping.

```
1 is_sleeping = atus_df.activity_name == 'Sleeping'
2 sleep_df = atus_df[is_sleeping]
3
4 avg_sleep = sleep_df.TUACTDUR.mean()
5
6 # even more info
7 summary_stats = sleep_df.TUACTDUR.describe()
8
9 # idxmax gives index with max value
10 max_row = sleep_df.TUACTDUR.idxmax()
11
12 # compare
13 sleep_df.loc[max_row]
14 atus_df.loc[max_row]
```

```
1 a = atus_data.activity_name.value_counts()
2 b = atus_data.activity_name.value_counts(normalize = True)
```

12.4 Applications

12.4.1 Interview Question

Start with data like the first table and create the second table.

```

1 # create original table
2 shares = pd.DataFrame(np.random.dirichlet([1,1,1],
3     size = 100),
4     columns = ['cycling', 'running', 'chess'])

```

	cycling	running	chess
0	0.55	0.32	0.13
1	0.47	0.03	0.50
2	0.31	0.43	0.26

	discipline1	discipline2	discipline3
0	0.55	0.32	0.13
1	0.50	0.47	0.03
2	0.43	0.31	0.26

Solution: It's homework!

This exercise highlights the difference between the structure of a DataFrame and a NumPy array. This task can be done most easily by converting the DataFrame to a NumPy array first. A more obvious solution involves looping through the DataFrame, row by row. Can you think of other ways to avoid loops? One route might use the fact that we have just three columns, so we are dealing with a minimum value, a maximum value, and then use the fact that the last value must add to the min and max to make one.

Chapter 13

Pandas: Join, Merge, and Other Manipulation

Reference: VanderPlas 2016a

13.1 Application: Primitive Pandas

This section doesn't exactly match anything in the book. We'll apply some of what we learned previously and consider some old-fashioned ways to loop through a DataFrame.

Consider the `rock_paper_scissors.csv` dataset. Load it using `pd.read_csv`. Each row represents a unique person and their strategy in a game of Rock, Paper, Scissors, where a strategy is just the chance they select either rock, paper, or scissors.

Let's verify that the probabilities sum to one.

1. Do this with a for loop.
2. Do this with the `sum` method.

Let's look for individuals who have very uneven strategies, in the sense that they lean toward any of the three actions with a chance greater than 0.5.

1. Do this with a for loop.
2. Do this with the `max` method.

Create three new columns that are agnostic of the specific rock, paper, or scissors actions and instead give the highest share, the second highest share, and the lowest share.

13.2 The Basic Join

We will cover merges and joins more in depth in the future, but for now let's consider the special case of joining two DataFrames.

DataFrames have a `join` instance for merging by the index. This requires similar indices and non-overlapping columns.

```
1 ser1 = pd.Series({'Alice':0.3, 'Bob':0.6})
2 ser2 = pd.Series({'Alice':0.7})
3
4 df1 = pd.DataFrame(ser1)
5 df2 = pd.DataFrame(ser2)
6
7 # This will fail
```

```

8 #joined = df1.join(df2)
9
10 # Rename columns
11 df1.columns = ['Rock']
12 df2.columns = ['Other']
13
14 joined = df1.join(df2)
15 joined2 = df1.join(df2, how = 'outer')
16
17 # examine
18 joined1 = df1.join(df2)
19 joined2 = df2.join(df1)
20 joined3 = df2.join(df1, how = 'outer')
21
22 print(joined1)
23 print(joined2)
24 print(joined3)
25
26 print(joined3.dropna())

```

13.3 Data Aggregation and Group Operations

13.3.1 GroupBy

The `groupby` method is fundamental to many data summary tasks. Load `purchase_transactions.csv`.¹ This dataset contains a row for each transaction, with `id` identifying the customer and `item` and `spent` giving the purchased item and spent giving the amount spent (corresponding to a quantity).

It'd be natural for us summarize this data by the individual customer. This requires creating a *GroupBy* object using the `groupby` method.

DataFrame Group By Object

```

1 grouped = df.groupby('id')
2 grouped.mean() # try this

```

Series Group By Object

```

1 grouped = df.groupby('id')['item']
2 grouped.mean() # try this

```

Group By With Multiple Columns

You can group across multiple dimensions by passing a list into the `groupby` method.

```

1 grouped2 = df.groupby(['id', 'item']) # DataFrame groupby object
2
3 grouped_df = grouped2.mean() # Creates a DataFrame
4
5 grouped_df # inspect

```

This kind of `groupby` creates a *MultiIndex*, even if you group a series instead of the whole DataFrame. Print `grouped_df.index` to see the following.

```

1 MultiIndex([( 0,  'apple'),
2             ( 1,  'butter'),
3             ( 2,  'apple'),
4             ( 2,  'butter'),

```

¹This is a randomly generated dataset that came from the [lifetimes](#) package and then I added extra columns. While more familiar topics might be easier, I would recommend this as a presentation subject for anyone interested in lifetime value calculations.


```

5      ( 2, 'orange'),
6      ( 2, 'turnip'),
7      ( 3, 'turnip'),
8      ( 4, 'orange'),
9      ( 4, 'turnip'),
10     ( 5, 'apple'),
11     ...
12     (4994, 'orange'),
13     (4994, 'turnip'),
14     (4995, 'orange'),
15     (4996, 'apple'),
16     (4996, 'butter'),
17     (4996, 'orange'),
18     (4996, 'turnip'),
19     (4997, 'turnip'),
20     (4998, 'apple'),
21     (4999, 'turnip')],
22     names=['id', 'item'], length=8986)

```

While the index of `df` was a list of integers, this is a list of tuples. A row of `grouped_df` is accessed with the standard loc, `grouped_df.loc[(0, 'apple')]`.

MultiIndices can complicate your code. You can get rid of the MultiIndex with `unstack`.

```
1 grouped_df.unstack()
```

Now, the values in the second dimension of the index become columns.

13.4 Concat and Append

The simplest way to combine two datasets is by concatenating them. Appending one dataset to another can be like a SQL union.

First, there is the `append` method. Consider the two American Time Use Survey datasets, `ATUS_activity_2018.csv` and `ATUS_activity_2019.csv`. It might be natural to combine these datasets if we don't see an important difference between 2018 and 2019. And even if there is an important difference, that could be noted by an extra column indicating the year.

With the `append` method², we can simply call `df2018.append(df2019)`. Note this returns a new DataFrame. You might assign this to a new variable if you'd like to work with the combined data.

```

1 df1819 = df2018.append(df2019)
2
3 # Compare number of rows
4 print(len(df2018), len(df2019))
5 print(len(df1819))
6
7 # Compare number of columns
8 print(len(df2018.columns), len(df2019.columns))
9 print(len(df1819.columns))

```

There is also the pandas function `concat`. We can use this to concatenate Series or DataFrames. The objects to be concatenated must be passed as a sequence and there is an optional `axis` argument.

```

1 # pd.concat(df2018, df2019) # Doesn't work
2
3 df_a = pd.concat([df2018, df2019]) # vertical
4 df_b = pd.concat([df2018, df2019], axis = 0) # vertical
5 df_c = pd.concat([df2018, df2019], axis = 1) # horizontal

```

Print out these DataFrames and compare the shapes. Then, inspect the indices. Note that for `df1819`, `df_a`, and `df_b`, the index now contains duplicates. You might amend this with `.reset_index()`. Pandas concatenation preserves indices. You can handle this by

1. Catching duplicates as an error

²See VanderPlas Chapter 3. Unlike the list `append`, this does not modify the original object.

2. Ignoring the Index
3. Adding MultiIndex keys.

To throw an error if there are duplicates, use the argument `verify_integrity = True`. To ignore the index, specify `ignore_index = True`. Or, to create a MultiIndex, pass an argument `keys = [2018, 2019]` where `keys` could more generally be any list that gives a unique key for each input to the concatenation.

Finally, there is also the `join` argument for `concat()` which can be used when the DataFrames don't share all of their columns. Use `join = 'inner'` to return just the common columns, and use `join = 'outer'` (also the default) to return all columns.

13.4.1 Application: What precedes sleeplessness?

Concatenation can be useful when you want to add columns that give values from the previous row. As an analyst, you might want to compare a row event with what took place before. We can do this with the help of the `shift` method.

```
1 df1819.reset_index(drop = True, inplace = True) # Clean index
2
3 shift_df1819 = df1819[['TUCASEID', 'activity_name']].shift() # pushes every row forward by
   default
4 shift_df1819.columns = ['prev_TUCASEID', 'prev_activity_name']
5
6 df1819 = pd.concat([df1819, shift_df1819], axis = 1)
7
8 df1819.head()
```

Then,

```
1 same = df1819.TUCASEID == df1819.prev_TUCASEID
2 sleepless = df1819.activity_name == 'Sleeplessness'
3 sleeping = df1819.activity_name == 'Sleeping'
```

Compare `df1819[same & sleepless].prev_activity_name.value_counts(normalize = True)` and `df1819[same & sleeping].prev_activity_name.value_counts(normalize = True)`.

13.5 Merge

We previously looked at the pandas join, which merged DataFrames based on their indices. Now, we will consider merges more generally, where we can merge based on column values.

A merge can be accomplished with a `.merge()` method, `df1.merge(df2 ...)` or with the pandas merge function, `pd.merge(df1, df2, ...)`.

First, let's consider `pd.merge` and let's load the 2018 ATUS data files. These DataFrames share just one column, `TUCASEID`.

```
1 activity = pd.read_csv("ATUS_activity_2018.csv", index_col = 'Unnamed: 0')
2 resp = pd.read_csv("ATUS_respondent_2018.csv", index_col = 'Unnamed: 0')
3
4 merge1 = pd.merge(activity, resp)
```

We didn't specify a column to merge on, so the merge is automatically done on the common column. However, it is better to specify using the `on` argument. This is to follow the principle of coding, "Explicitly is better than implicit."

As we could use the `verify_integrity` argument in concatenation, we can use a `validate` argument to throw an error if Python doesn't find our expected behavior in the merge. Here, we have a many-to-one merge, because a single `TUCASEID` appears multiple times in the left dataset, `activity`, and just once in the right dataset, `resp`.

```
1 try:
2     pd.merge(activity, resp, on = 'TUCASEID', validate = 'many_to_one')
3 except Exception as e:
4     print(e)
```

```

5
6 try:
7     pd.merge(activity, resp, on = 'TUCASEID', validate = 'one_to_many')
8 except Exception as e:
9     print(e)
10
11 try:
12     pd.merge(activity, resp, on = 'TUCASEID', validate = 'one_to_one')
13 except Exception as e:
14     print(e)

```

Merges can also be done on multiple columns. Here we create (fake) supplemental data to be added.

```

1 activity_supplement = activity[['TUCASEID', 'TUSTARTTIM']]
2 activity_supplement.loc[:, 'is_alone'] = np.random.choice([True, False], len(
    activity_supplement))
3
4 # Jumble dataframe to a simple index join or concat is possible
5 activity_supplement = activity_supplement.sample(len(activity_supplement)) # samples without
    replacement to shuffle
6 activity_supplement.reset_index(drop = True, inplace = True)
7
8 # Merge
9 pd.merge(activity, activity_supplement, on = ['TUCASEID', 'TUSTARTTIM'], validate = '
    one_to_one')

```

Datasets can also be merged on differently named columns.

```

1 activity_supplement.columns = ['a', 'b', 'c']
2 pd.merge(activity, activity_supplement, left_on = ['TUCASEID', 'TUSTARTTIM'], right_on = ['a
    ', 'b'])

```

Finally, there are left, right, inner, and outer joins. These can be specified with `how`.

Consider the following

```

1 pd.merge(activity, activity_supplement.head(), left_on = ['TUCASEID', 'TUSTARTTIM'],
    right_on = ['a', 'b'])

```

What results is a DataFrame of length. By default, pandas does an inner join. Inspect each DataFrame below. These are all unique!

```

1 supplement2 = activity_supplement.head()
2 supplement2.columns = ['TUCASEID', 'TUSTARTTIM', 'is_alone'] # name back
3
4 # Create a TUCASEID not in the activity dataset
5 supplement2.loc[0, 'TUCASEID'] = "Uncle Milton"
6
7
8 df_inner = pd.merge(activity, supplement2, on = ['TUCASEID', 'TUSTARTTIM'], how = 'inner')
9
10 df_outer = pd.merge(activity, supplement2, on = ['TUCASEID', 'TUSTARTTIM'], how = 'outer')
11
12 df_right = pd.merge(activity, supplement2, on = ['TUCASEID', 'TUSTARTTIM'], how = 'right')
13
14 df_left = pd.merge(activity, supplement2, on = ['TUCASEID', 'TUSTARTTIM'], how = 'left')

```

Note that in cases where they may be no conceptual difference between a left and right merge, you might find the rows are still ordered differently so identical DataFrames are not returned. Below, see the right (left) merge creates a new DataFrame with the same order as in the right (left) DataFrame. This is also our first use of the DataFrame method `equals()`, which tests for the equality of all elements.

```

1 left = pd.DataFrame({'id': [1,2],
2                     'country_name': ['USA', 'CN']})
3 right = pd.DataFrame({'gold_medals': [9,8],
4                      'id': [2,1]})
5
6 left_merged = left.merge(right, on = 'id', how = 'left')
7 right_merged = left.merge(right, on = 'id', how = 'right')
8

```

```

9 # This returns False
10 print(left_merged.equals(right_merged))
11
12 # This returns True
13 print(left_merged.sort_values('id').reset_index(drop = True).equals(\
14     right_merged.sort_values('id').reset_index(drop = True)))

```

13.6 Pivot Tables

VanderPlas 2016a motivates pivot tables as a kind of advanced `groupBy` operation. Below, we slightly modify the example on page 171.

```

1 import seaborn as sns
2 titanic = sns.load_dataset('titanic')
3 titanic.groupby(['sex', 'class']).survived.mean()

```

This creates a series with a MultiIndex. You might prefer to have the same data in a table. This can be done by adding `.unstack()`. That is, run `titanic.groupby(['sex', 'class']).survived.mean().unstack()`.

`pivot_table()`

```

1 # DataFrame method
2 p1 = titanic.pivot_table('survived', index = 'sex', columns = 'class')
3
4 # Equivalent pandas function
5 p2 = pd.pivot_table(titanic, values = 'survived', index = 'sex', columns = 'class', aggfunc
6     = np.mean)

```

`pivot()`

Another method is `df.pivot()`. This can't handle duplicates. It doesn't do any aggregation. The data is simply reshaped.

13.6.1 Crosstabs

The `crosstab()` function simplifies a pivot with `aggfunc = 'count'`.

```

1 df = pd.read_csv("purchase_transactions.csv", index_col = 'Unnamed: 0')
2
3 # pivot
4 pd.pivot_table(df, index = 'id', columns = 'item', values = 'spent',
5     aggfunc='count', fill_value = 0).head()
6
7 # equivalent
8 pd.crosstab(df1.id, df1.item)

```

Chapter 14

Data Visualization

There are many ways to skin a cat and there are many ways to plot data in Python. But using [matplotlib](#) is the obvious route and what we'll focus on. Matplotlib offers a major advantage of being integrated into pandas and being the most popular choice for anyone working with the broader Python community. [Seaborn](#) is built on matplotlib and might be prettier out of the box. [Plotnine](#) is based off ggplot2, so it might be more comfortable turf for R users.

While you're getting started, you'll likely be tempted to use MS Excel or Google Sheets instead of using matplotlib. Try to fight through. In the long-run, matplotlib will be more re-usable and provide more customization.

Our reference for matplotlib is VanderPlas [2016a](#) Chapter 4. For another tutorial, check out [Nicolas Rougier's tutorial](#). For book-length treatments, see Rougier [2021](#) or Clark [2022](#). I have a few random tutorials on [my youtube channel](#).

There are two interfaces: a MATLAB-inspired interface and an object-oriented interface. That is, you can create plots with either of this code styles. First, we will work with the simpler MATLAB style.

14.1 MATLAB Interface

The standard import and alias is as follows.

```
1 import matplotlib.pyplot as plt
```

Here's a minimal working example for a plot.

```
1 import numpy as np
2
3 x = np.linspace(0, 10, 100)
4 plt.plot(x, np.sin(x))
5 plt.plot(x, np.cos(x))
6
7 plt.show()
```

The `plot()` function makes a line plot. Using `plt.show()` is not always necessary in a Jupyter environment. See [VanderPlas Chapter 4](#) for a slightly more technical discussion. It's sufficient to think of `plt.show()` as “finishing” the plot. This can be helpful when you want to create multiple plots in one code block.

Compare the following two programs.

```
1 x = np.linspace(0, 10, 100)
2
3 for i in range(0, 10):
4
5     y = np.ones(len(x)) * i
6     plt.plot(x, y)
7     plt.show()
```

```
1 x = np.linspace(0, 10, 100)
2
```

```

3 for i in range(0,10):
4
5     y = np.ones(len(x)) * i
6     plt.plot(x, y)
7
8 plt.show()

```

The only difference is the indentation of `plt.show()`. Observe the first creates ten different plots. The second creates just one plot.

14.1.1 Saving Figures

You can use `plt.savefig` to save a plot. This must go before `plt.show()`. There are many supported filetypes. I am partial to png and pdf. A `dpi` argument can be used to set the image resolution.

```

1 x = np.linspace(0,10,100)
2
3 for i in range(0,10):
4
5     y = np.ones(len(x)) * i
6     plt.plot(x, y)
7
8 plt.savefig('example_figure.png', dpi = 300)
9 plt.show()

```

14.1.2 Special Plots: Scatter, Histogram, and Bar Plots

Scatter plots can be created with `plt.scatter()`. However, you can get some efficiency gain by instead using `plt.plot()` by specifying a marker, like `marker = 'o'`, and setting `linestyle = ''`. See VanderPlas 2016a for further discussion.

Plot type	Function
Scatter	<code>plt.scatter(x,y, alpha = 0.1)</code>
Histogram	<code>plt.hist(x, bins = 20)</code>
Bar Plot	<code>plt.bar(x,y)</code>
Horizontal Bar Plot	<code>plt.barh(x,y)</code>

14.1.3 Customizations

There are many additional parameters for customization available in each of these plotting functions and methods. It would be a fool's errand to learn them all from the outset. Over time you will find yourself visiting matplotlib documentation and getting a sense of what's possible and what's not.

You should be familiar with

Function, Method, or Keyword Parameter	What it does	Example Availability
<code>alpha</code>	set opacity	<code>plt.scatter()</code>
<code>figsize</code>	figure dimensions	<code>plt.subplots()</code> <code>plt.figure()</code>
<code>dpi</code>	figure resolution	<code>plt.savefig()</code>
<code>transparent</code>	transparent background	<code>plt.savefig()</code>
-	set aspect ratio	<code>ax.set_aspect()</code>
-	vertical line	<code>ax.axvline</code> <code>plt.axvline()</code>
-	horizontal line	<code>ax.axhline</code> <code>plt.axhline()</code>

14.1.4 Pandas Integration

There are Pandas methods that create Matplotlib objects. You'll get pretty far with `.plot()` and `.plot.bar()` or `.plot.barh()`.

Let's return to our sleeplessness application. We can call `.plot.bar()` for a quick bar graph.

```

1 top_activities_before_sleeping = df1819[same & sleeping].prev_activity_name.value_counts(
    normalize = True).head()
2 top_activities_before_sleepless = df1819[same & sleepless].prev_activity_name.value_counts(
    normalize = True).head()
3
4 top_activities_before_sleeping.plot.bar()

```

These don't quite like in the first example where we could keep adding to a plot. For example, the below does not create a side-by-side barplot or even create accurate labels on the x-axis to account for the differing indices.

```

1 top_activities_before_sleeping = df1819[same & sleeping].prev_activity_name.value_counts(
    normalize = True).head()
2 top_activities_before_sleepless = df1819[same & sleepless].prev_activity_name.value_counts(
    normalize = True).head()
3
4 # This is bad!
5 top_activities_before_sleeping.plot.bar(color = 'red', alpha = 0.5)
6 top_activities_before_sleepless.plot.bar(color = 'blue', alpha = 0.5)

```

For a side-by-side plot, it is better to use a DataFrame instead of a Series.

```

1 a = pd.DataFrame(top_activities_before_sleeping)
2 a.columns = ['sleeping']
3
4 b = pd.DataFrame(top_activities_before_sleepless)
5 b.columns = ['sleepless']
6
7 a.join(b, how = 'inner').plot.bar()

```

14.2 Object-oriented Interface

14.2.1 Introduction to the OO Interface

This section is from Clark 2022.
[Github Link](#)

The object-oriented interface looks like this.

```

1 fig, ax = plt.figure(), plt.axes()
2 ax.plot(x,y)
3 ax.set_title("My Chart")

```

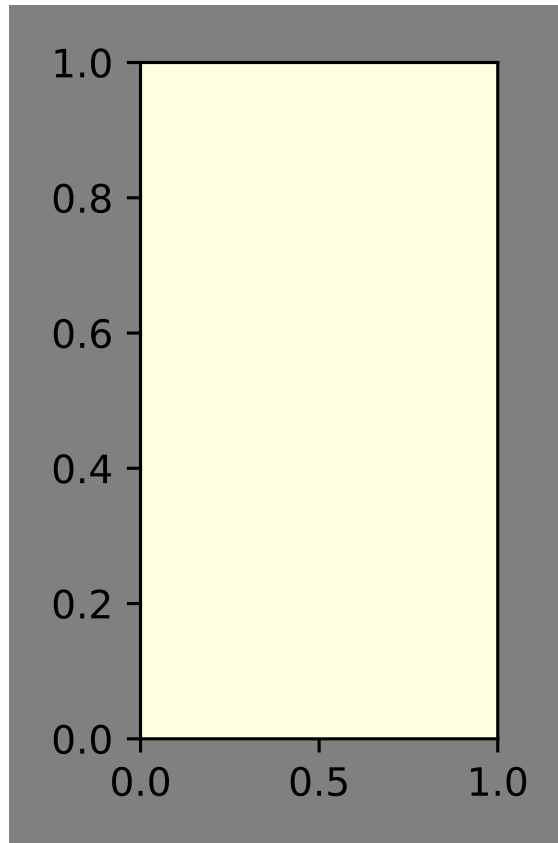
There is no such thing as a free lunch, so you will observe this interface requires more code to do the same exact thing. Its virtues will be more apparent later. Object-oriented programming (OOP) also requires some new vocabulary. OOP might be contrasted with procedural programming as another common method of programming. In procedural programming, the MATLAB-style interface being an example, the data and code are separate and the programmer creates procedures that operate on the program's data. OOP instead focuses on the creation of *objects* which encapsulate both data and procedures.

An object's data are called its *attributes* and the procedures or functions are called *methods*. In the previous code, we have figure and axes objects, making use of axes methods `plot()` and `set_title()`, both of which add data to the axes object in some sense, as we could extract the lines and title from `ax` with more code. Objects themselves are instances of a *class*. So `ax` is an object and an instance of the Axes class. Classes can also branch into subclasses, meaning a particular kind of object might also belong to a more general class. A deeper knowledge is beyond our scope, but this establishes enough vocabulary for us to continue building an applied knowledge of matplotlib. Because `ax` contains its data, you can think of `set_title()` as changing `ax` and this helps make sense of the `get_title()` method, which simply returns the title belonging to `ax`. Having some understanding that these objects contain both procedures and data will be helpful in starting to make sense of intimidating programs or inscrutable documentation you might come across.

Figure, Axes

A plot requires a figure object and an axes object, typically defined as `fig` and `ax`. The figure object is the top level container. In many cases like in the above, you'll define it at the beginning of your code and never need to reference it again, as plotting is usually done with axes methods. A commonly used figure parameter is `figsize`, to which you can pass a sequence to alter the size of the figure. Both the figure and axes objects have a `facecolor` parameter which might help to illustrate the difference between the axes and figure.

```
1 fig = plt.figure(figsize = (2,3),
2                 facecolor = 'gray')
3 ax = plt.axes(facecolor = 'lightyellow')
```



The axes object, named `ax` by convention, gets more use in most programs. In place of `plt.plot()`, you'll use `ax.plot()`. Similarly, `plt.hist()` is replaced with `ax.hist()` to create a histogram. If you have experience with the MATLAB interface, you might get reasonably far with the object-oriented style just replacing the `plt` prefix on your pyplot functions with `ax` to see if you have an equivalent axes method.

This wishful coding won't take you everywhere though. For example, `plt.xlim()` is replaced by `ax.set_xlim()` to set the *x*-axis view limits. To modify the title, `plt.title()` is replaced with `ax.set_title()` and there is `ax.get_title()` simply to get the title. The axes object also happens to have a `title` attribute, which is only used to access the title, similar to the `get_title()` method. Many matplotlib methods can be classified as *getters* or *setters* like for these title methods. The plot method and its logic is different. Later calls of `ax.plot()` don't overwrite earlier calls and there is not the same getter and setter form. There's a `plot()` method but no single `plot` attribute being mutated. Whatever has been plotted can be retrieved, or gotten (getter'd?), but it's more complicated and rarely necessary. Use the code below to see what happens with two calls of `plot()` and two calls of `set_title()`. The second print statement demonstrates that the second call of `set_title()` overwrites the title attribute, but a second plot does not nullify the first.

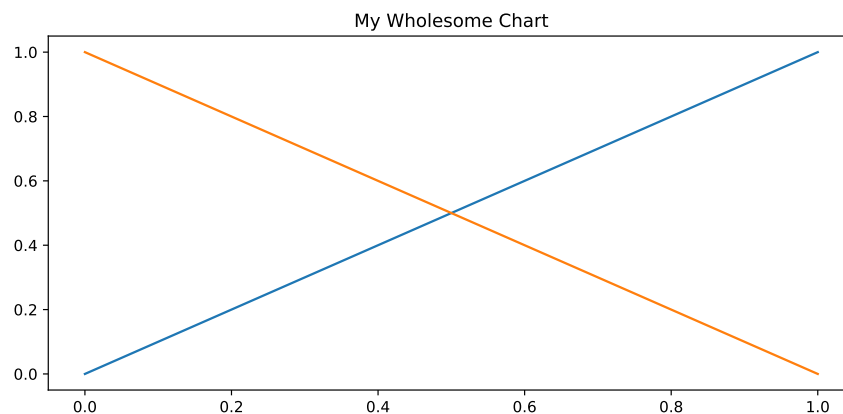
```
1 x = np.linspace(0,1,2)
2 fig, ax = plt.figure(figsize = (8,4)), plt.axes()
```



```

3 ax.plot(x, x)
4 ax.plot(x, 1 - x)
5 ax.set_title("My Chart")
6 print(ax.title)
7 print(ax.get_title()) # Similar to above line
8 ax.set_title("My Wholesome Chart")
9 print(ax.get_title()) # long

```



Axes methods `set_xlim()` and `get_xlim()` behave just like `set_title()` and `get_title()`, but note there is no attribute simply accessible with `ax.xlim`, so the existence of getters and setters is the more fundamental pattern.¹

Mixing the Interfaces

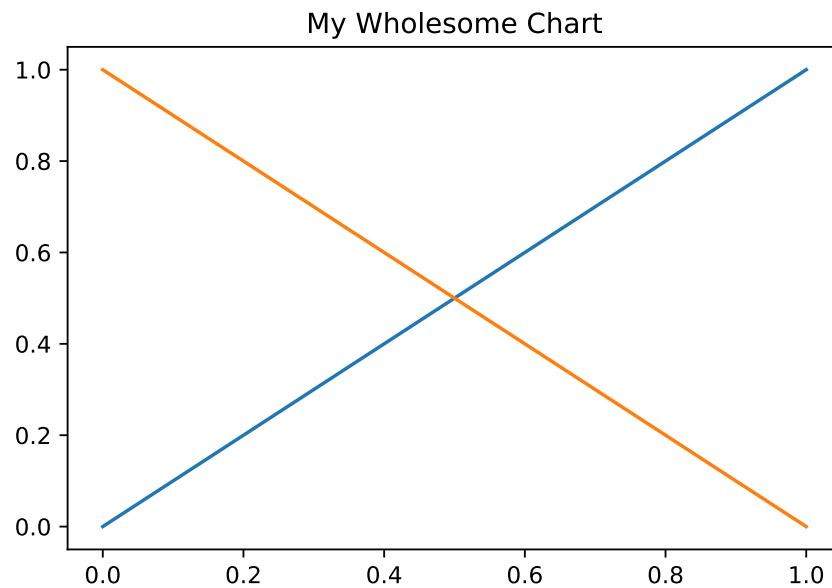
You can also mix the interfaces. Use `plt.gca()` to get the current *axis*. Use `plt.gcf()` to get the current *figure*.

```

1 x = np.linspace(0,1,2)
2 plt.plot(x,x)
3 plt.title("My Chart")
4
5 ax = plt.gca()
6 print(ax.title)
7
8 ax.plot(x, 1 - x)
9 ax.set_title('My Wholesome Chart')
10 print(ax.title)
11
12 fig = plt.gcf()
13 fig.savefig('chart.pdf') # same as plt.savefig

```

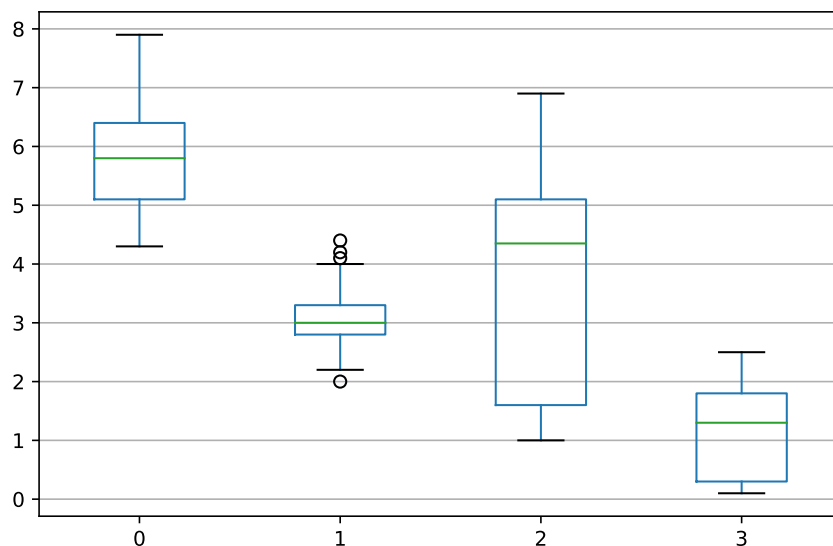
¹Getters and setters are thought of as old-fashioned. It's more Pythonic to access attributes directly, but matplotlib doesn't yet support this.



In the above, we started with MATLAB and then converted to object-oriented. We can also go in the opposite direction as well. We can also mix pandas with MATLAB or OOP-style matplotlib.

These plots can be mixed with the object-oriented interface. You can use a plot method and specify the appropriate axes object as an argument. Below we import the iris dataset and make a boxplot with a mix of axes methods and then pyplot functions.

```
1 from sklearn.datasets import load_iris
2 data = load_iris()['data']
3 df = pd.DataFrame(data)
4
5 fig, ax = plt.figure(), plt.axes()
6
7 df.plot.box(ax = ax)
8 ax.yaxis.grid(True)
9 ax.xaxis.grid(False)
10
11 plt.tight_layout()
12 plt.savefig('irisbox.pdf')
```



14.2.2 Subplots with `plt.subplots()`

You can use `subplots()` to create multiple subplots within the figure.

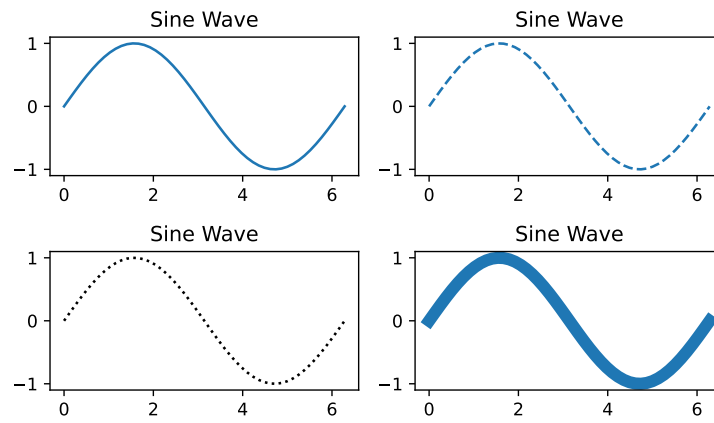
```
1 # Multiple Subplots
2 # ax is a tuple for two different axes
3 fig, ax = plt.subplots(1,2)
4
5 # Call plot() on the axis
6 ax[0].plot(x, np.sin(x))
7 ax[1].plot(x, np.sin(x), color = 'tomato')
8
9 plt.show()
```

In the last example above, two subplots are created. If you save this, note this creates a single image file, not one for each subplot. Next, we create a 2×2 grid. Note, when we do this, we must index `ax` at two layers. Before, in a 1×2 , we had one-dimensional indexing. The superfluous second dimension was *squeezed* out by default.²

```
1 fig, ax = plt.subplots(2,2)
2
3 x = np.linspace(0, 2 * np.pi)
4 y = np.sin(x)
5
6 ax[0,0].plot(x,y)
7 ax[0,1].plot(x,y, linestyle = 'dashed')
8 ax[1,0].plot(x,y, linestyle = 'dotted', color = 'black')
9 ax[1,1].plot(x,y, linewidth = 7)
10
11 for ax_ in fig.axes:
12     ax_.set_title('Sine Wave')
13
14 fig.suptitle("Big Plot", size = 20)
15 plt.tight_layout()
16
17 plt.savefig("subplot_example_2by2.pdf")
```

²See the `squeeze` parameter in the [documentation](#).

Big Plot



Chapter 15

Time and Dates

Additional reference: VanderPlas 2016a Chapter 3

Up to now, we’ve only encountered dates as strings like `independence_day = '1776-07-04'`. This is limited. For example, there are no string operations to add and subtract days. To that end, it is helpful to have new classes for date objects and ways to manipulate those dates, times, and datetimes. The built-in `datetime` and `dateutil` modules are useful for that. These will only support times since 1 AD. This might disappoint astronomers or anyone else wanting to extend their analysis to years BC. If that’s you, you’ll have to study the AstroPy [Time class](#).

15.1 Using the datetime module

The standard import is `import datetime as dt`.

15.1.1 Dates and Datetimes

The two main classes are `date` and `datetime`. A `date` object is created with `dt.date(year, month, day)`. `datetime` extends this with optional parameters for hour, minute, second, microsecond, timezone information, and [fold](#).

```
1 last_lecture1 = dt.date(2022, 5, 2)
2 last_lecture2 = dt.datetime(2022, 5, 2, 20, 10)
3
4 now = dt.datetime.today() # get current date and time
5 print(now) # local NYC time for me
```

You can go back to a date from a datetime with the datetime method `date()`. You can get a `time` object with the datetime method `time()`.

We won’t cover time zones or folds. But on time zones, you should be aware that some datetimes are automatically in local time. In New York City, this comes as a -4 or -5 offset from UTC (coordinated universal time). For more with timezones, use [pytz](#). Below is a small taste of `pytz` for anyone curious.

```
1 import pytz
2
3 unaware = dt.datetime(2014,2,2) # time zone unaware datetime
4 aware = pytz.utc.localize(unaware) # sets tzinfo to be UTC without changing the hour
```

15.1.2 Unix Timestamps

Beyond strings like `'2022-12-25'`, you might also encounter timestamps as integers. That is, you might use [Unix time](#). This is also called POSIX time. These integers give the number of seconds elapsed since 00:00:00 UTC on 1 January 1970 (the Unix epoch).

These can be converted to datetime objects as shown below.

```

1 # UTC
2 # This gives 1970, 1, 1, 0, 0
3 dt.datetime.utcnow().timestamp(0)
4
5 # Local Time
6 # This gives 1969, 12, 31, 19, 0 (NYC)
7 dt.datetime.fromtimestamp(0)
8
9 # Can also use dt.date on fromtimestamp
10 # This gives 1969, 12, 31 (NYC)
11 dt.date.fromtimestamp(0)

```

15.1.3 Time Deltas and Relative Deltas

A `timedelta` object is useful for dealing with a duration of time, when subtracting datetimes for example.

```

1 class_start = dt.datetime(2020,12,5,12,10,0)
2 class_end = class_start + dt.timedelta(hours = 1, minutes = 50)
3
4 class_duration = class_end - class_start
5
6 print(class_duration)

```

The `seconds` attribute is different than the `total_seconds()` method. Can you anticipate the difference?

```

1 delta = dt.timedelta(days = 1, seconds = 100)
2
3 print(delta.seconds)
4 print(delta.total_seconds())

```

But what if we just want the next month? We can't use a single `timedelta` that takes us from February 1st to March 1st and also from April 1st to May 1st. For this, we need a `relativedelta` from `dateutil`.

```

1 from dateutil.relativedelta import relativedelta
2
3 date = dt.date(2019, 12, 9)
4 other_date = date + relativedelta(months = 2)
5 print(other_date)
6
7 # Consider how this works with leap year
8 next_year1 = date + relativedelta(years = 1)
9 next_year2 = date + dt.timedelta(days = 365)

```

15.1.4 Date Strings

The datetime `strptime` method is very useful in string conversion. It works with a date string as its first argument and a `format` as a second argument.

For example, `dt.datetime.strptime('2020-01-01', '%Y-%m-%d')` returns a datetime object for January 1st, 2020.

Above, the string `'%Y-%m-%d'` is a date format code. Here are some common format codes, applied to Sunday January 30, 2000, 11:59PM, local to Louisville, Kentucky. These can all be verified with `pd.Timestamp` (`year = 2000, month = 1, day = 30, hour = 23, minute = 59, tz = 'America/Kentucky/Louisville'`).`strptime()`.

Code	Output/Example
'%Y'	4-Digit Year
'%m'	Month Number
'%d'	Day of Month
'%B'	Month Name
'%H'	24-Hour Clock Hour
'%M'	Minute
'%h'	12-Hour Clock Hour
'%p'	AM or PM
'%A'	Day of Week
'%Z'	Timezone Name
'%Y-%m'	'2000-01'
'%Y/%m/%d'	'2000/01/30'
'%B %y'	'January 00'
'%H:%M %Z'	'23:59 EST'
'%A %I%p'	'Sunday 11PM'

A more complete list of format codes can be found at strftime.org. Codes that generate actual names, like '%A' or '%B', can be made lowercase to produce an abbreviated name. Notice that these formats create zero-padded numbers like '07' instead of '7'. On Mac or Linux, padding can be eliminated with the '-' modifier, using '%-H' or '%-m' instead of '%H' or '%m' for example. On Windows, use '#'.

Dateutil

Dateutil offers, among other things, a details-free parser (you don't have to specify a format code).

```
1 from dateutil import parser
2 date = parser.parse("9th of December, 2019")
3 print(date)
```

15.1.5 Pandas and Numpy

Pandas offers an efficient `Timestamp` object and NumPy offers an efficient `datetime64`. For NumPy, the tradeoff is they are less flexible than datetime objects. Pandas offers something closer to the best of both worlds.

Notice the accessing a NumPy array from a Series will convert datetime objects into NumPy's `datetime64`.

```
1 ser = pd.Series([dt.datetime.today()])
2 date = ser.values[0]
3 print(date)
4
5 # date + dt.timedelta(days = 1) # error
```

Here, we parse a date into a Pandas Timestamp and use a datetime timedelta.

```
1 date = pd.to_datetime("4th of July, 2015")
2 print(date)
3 print(type(date))
4
5 date + dt.timedelta(days = 1) # Works
```

In pandas, we can also parse dates when reading in a CSV, using the `parse_dates` parameter in `pd.read_csv()`. A string for a single column or a list for multiple columns are both valid arguments.

15.2 Analysis

```
1 dates = pd.date_range(start='2020-12-05', end='2021-12-05', freq='1D')
2
3 df = pd.DataFrame(index = dates)
4
5 df['stock_price'] = range(len(dates)) + np.random.normal(0,10,len(dates))
```

```

6 plt.plot(df['date'], df['stock_price'])
7 plt.show()
8
9
10 ## Rolling Function
11
12 data = df.stock_price.rolling(30, center = True)
13 data.mean().plot()
14 plt.show()

```

Try it with other window sizes.

```

1 data = df.stock_price.rolling(365, center = True)
2 data.mean().plot()
3 plt.show()
4
5 data = df.stock_price.rolling(100, center = True)
6 data.mean().plot()
7 plt.show()

```

15.2.1 Application 1

Load the NYC taxi data and look at the distribution of ride lengths.

```

1 df = pd.read_csv('nyc_data.csv')
2
3 df.head(1)
4
5
6 df['pickup_datetime'] = pd.to_datetime(df['pickup_datetime'])
7 df['dropoff_datetime'] = pd.to_datetime(df['dropoff_datetime'])
8
9 df['duration'] = df['dropoff_datetime'] - df['pickup_datetime']
10
11 # Make histogram with logarithmic y axis
12 df['duration'].dt.total_seconds().hist(log = True)

```

15.2.2 Application 2

Load the ATUS_activity_2019.csv dataset. Create a new datetime time object column from the TUSTARTTIM column.

Method 1

```

1 import datetime
2 import pandas as pd
3
4 df = pd.read_csv("ATUS_activity_2019.csv")
5
6 # Parse the time
7 df['time_col'] = df['TUSTARTTIM'].map(lambda x: dt.datetime.strptime(x, "%H:%M:%S"))
8
9 # This sets everything to a datetime object from January 1900
10 # Convert to time with the time method
11
12 df['time_col'] = df.time_col.map(lambda x: x.time())

```


Chapter 16

Python for Excel

Reference: Zumstein 2021

Excel (or Google Sheets) is unavoidable if you work with enough people, especially ThinkPad people. You will want to deliver a spreadsheet to a stakeholder, not a Python script. And that spreadsheet will be received better if it's well formatted. That means we're not talking about a CSV file fresh from `df.to_csv()`. In this chapter, we'll use Python to do the following.

1. Save to individual sheets in a larger workbook.
2. Read in individual sheets as dataframes.
3. Apply cell formatting.
4. Apply filters.
5. Insert plots in a sheet.

If you're in this class, you're way too smart to do these tasks by pointing and clicking in a GUI. We'll use pandas and [OpenPyXL](#).

16.1 Pandas

Pandas comes with the `read_excel()` function and the `to_excel()` method. Below, we use these most basically.

```
1 df = pd.DataFrame({'a': [0]})
2
3 # Save to excel file
4 df.to_excel("mwe.xlsx", index = None)
5
6 # Clear df from memory
7 del df
8
9 # Read
10 df = pd.read_excel("mwe.xlsx")
```

We go into more detail in this section. However, pandas is limited. Separate reader and writer packages are useful for including charts, changing formatting, etc.

16.1.1 Reading

Above, `mwe.xlsx` is the simplest and most well-behaved kind of Excel file you can encounter. There are not multiple sheets and the data begins in cell A1. Often, you'll receive an Excel file with multiple sheets and with some formatting that, to a Python user, is wonky. Reading in the right sheets is priority number one. You can use the `sheet_name` parameter to specify the sheet you want to use, by name or index. If you pass a

list of sheet names, `read_excel()` returns a dictionary with the DataFrames as values and sheet names as the keys. Use `sheet_name = None` to load all of the sheets.

Maybe the first row and first column are blank, in which case the data begins in cell B2. Use `skiprows` and `usecols` to specify the cell range you want to use. These and other parameters are described in Table 7-1 in Zumstein 2021. The additional parameters are mostly useful for data cleaning, so you can have that done in the `read_excel()` call instead of in additional lines of code.

16.1.2 Writing

As with `read_excel()`, `to_excel()` includes a `sheet_name` parameter. However, this can't be used to add multiple sheets to a single Excel file. The method doesn't *add* a new sheet to an existing file, it merely writes the Excel file from scratch, with a single sheet named according to `sheet_name`.

```
1 df = pd.DataFrame({'a': [0]})
2
3 # creates mwe.xlsx with one sheet called A
4 df.to_excel("mwe.xlsx", index = None, sheet_name = 'A')
5
6 # overwrites mwe.xlsx, one sheet called B
7 df.to_excel("mwe.xlsx", index = None, sheet_name = 'B')
```

To add multiple DataFrames to different sheets (or even within the same sheet), you need to use the `ExcelWriter` class. Below, we use `ExcelWriter` as a *context manager*, using a `with` statement.

```
1 df_a = pd.DataFrame({'a': [0]})
2 df_b = pd.DataFrame({'b': [1]})
3
4 with pd.ExcelWriter('multi_sheet.xlsx') as writer:
5     df_a.to_excel(writer, sheet_name = 'A', index = None)
6     df_b.to_excel(writer, sheet_name = 'B', index = None)
```

The additional parameters for `to_excel()` deserve some attention, because they aren't as easily replaced by other lines of code. Neglecting them creates more work for you in Excel, and our objective is to avoid that. Some parameters are described in Table 7-2 in Zumstein 2021 and in the [official documentation](#).

Parameter	Description
<code>startrow</code>	First row where the DataFrame is written (using zero-based indexing)
<code>startcol</code>	First column where the DataFrame is written (using zero-based indexing)
<code>freeze_panes</code>	Takes a tuple for the number of rows and columns to freeze. Passing (1,2) freezes the first row and the first two columns.

16.2 OpenPyXL

Installation: OpenPyXL is useful for reading, writing, and editing Excel files. If you need to install it, run `conda install -c anaconda openpyxl` or `pip install openpyxl` in the terminal.

Remember, there are three levels to an Excel file: the workbook, the worksheet, and the individual cells. The workbook is the entire file and all the sheets (or tabs). The worksheet is an individual tab. And a cell is the entry at a specific coordinate in a worksheet. OpenPyXL has classes `Workbook`, `Worksheet`, and `Cell` for each of these.

16.2.1 Reading

Use `load_workbook()` to read in data and specify `data_only = True` to read in the cell values instead of the cell formulas.

```
1 import pandas as pd
2 import openpyxl
3 import datetime as dt
4
5 # Workbook object
```

```

6 book = openpyxl.load_workbook('multi_sheet.xlsx', data_only = True)
7
8 # get a list of all sheet names
9 print(book.sheetnames)
10
11 # Worksheet objects
12 sheet_a = book['A']
13 sheet_b = book['B']
14
15 # Get dimensions
16 print(sheet_a.max_row, sheet_a.max_column)
17
18 # Cell object
19 cell = sheet_a['A1']
20 print(cell.value)

```

16.2.2 Writing

We write a DataFrame to an Excel .xlsx file using the functionality provided directly in OpenPyXL, making use of `dataframe_to_rows`.¹ The process is as follows.

1. Create a Workbook instance.
2. Access a Worksheet. Below we use the `active` attribute to get the currently active sheet.
3. Append data to the bottom of the sheet using the Worksheet `append()` method.
4. Save the Workbook using the `save()` method.

```

1 from openpyxl.utils.dataframe import dataframe_to_rows
2
3 wb = openpyxl.Workbook()
4 ws = wb.active
5
6 for r in dataframe_to_rows(df, index=True, header=True):
7     ws.append(r)
8
9 wb.save("pandas_openpyxl.xlsx")

```

This includes the index and header (column name) values given the parameters in `dataframe_to_rows`. Our Excel file includes a sheet as shown below.

	A	B	C	D	E
1		Alice	Bob	Cathy	Dale
2					
3	Alice	1	0.793637	0.580004	0.162299
4	Bob	0.793637	1	0.500008	0.88952
5	Cathy	0.580004	0.500008	1	0.436747
6	Dale	0.162299	0.88952	0.436747	1

If we had simply used `df.to_excel('filename.xlsx')`, we'd end up with the below.

¹Zumstein 2021 provides a module [excel.py](#). It's handy for going back and forth between OpenPyXL and other Excel packages. However, it doesn't play that well with index and header values.

	A	B	C	D	E
1		Alice	Bob	Cathy	Dale
2	Alice	1	0.793637	0.580004	0.162299
3	Bob	0.793637	1	0.500008	0.88952
4	Cathy	0.580004	0.500008	1	0.436747
5	Dale	0.162299	0.88952	0.436747	1

Why the empty row when using OpenPyXL? I'm not sure.² The `dataframe_to_rows()` function finds an empty row somewhere, somehow. We can avoid appending this to our worksheet by modifying the for-loop, as is done below.

```

1 from openpyxl.utils.dataframe import dataframe_to_rows
2
3 wb = openpyxl.Workbook()
4 ws = wb.active
5
6 for r in dataframe_to_rows(df, index=True, header=True):
7     print(r)
8     if r != [None]:
9         ws.append(r)
10
11 wb.save("pandas_openpyxl.xlsx")

```

16.2.3 Styles

Styles and formatting can be applied, but only one cell at a time. Accordingly, we need to learn more about the `Cell` class.³ A specific cell can be accessed by indexing a worksheet with the appropriate coordinate, `ws['A1']` for example. Cell ranges can be accessed as might be done in Excel. Columns A-C can be obtained with `ws['A:C']`. Rows 3-5 can be accessed with `ws[3:5]`. Note cell ranges are just tuples of cells. There is also a `iter_rows()` worksheet method, returning a generator for iteration.

Below are some basic `Cell` attributes.

Attribute	Description
<code>coordinate</code>	The Excel coordinate of the cell (e.g. <code>'B2'</code>)
<code>column</code>	Column number of the cell (one-based)
<code>row</code>	Row number of the cell (one-based)
<code>column_letter</code>	Column letter of the cell (e.g. <code>'B'</code> or <code>'AA'</code>)
<code>value</code>	Data value of the cell
<code>is_date</code>	Boolean value

There are `Cell` style attributes you might use to inspect and to overwrite the cell font, border, alignment, fill, and number format. To set these, use the `openpyxl.styles` submodules. For example, to set the cell alignment, write `ws['C8'].alignment = openpyxl.styles.alignment.Alignment(horizontal = 'center')`.

Attribute	openpyxl.styles submodule and class
<code>font</code>	Font
<code>border</code>	<code>borders.Side</code> and <code>borders.Border</code>
<code>alignment</code>	<code>alignment.Alignment</code>
<code>fill</code>	<code>fills.PatternFill</code>
<code>number_format</code>	N/A (use a string)

²See the [documentation here](#) if you'd like to delve further on your own.

³[Cell documentation here](#).

The final attribute listed above, `number_format`, is set with a string. Available format examples can be found in the [OpenPyXL documentation](#) and the [Microsoft documentation](#) provides further detail. Below are a few formats and their Excel output.⁴

Format String	Example	Example (negative)
General	1984.1984	-1984.1984
0	1984	-1984
0.00	1984.20	-1984.20
#,##0	1,984	-1,984
#,##0.00	1,984.20	-1,984.20
+"\$"#,##0_);("\$"#,##0)	+\$1,984	(\$1,984)
"\$"#,##0_);[Red]("\$"#,##0)	\$1,984	(\$1,984)
[Green]"\$"#,##0_);[Red]("\$"#,##0)	\$1,984	(\$1,984)
"\$"#,##0.00_);("\$"#,##0.00)	\$1,984.20	(\$1,984.20)
"\$"#,##0.00_);[Red]("\$"#,##0.00)	\$1,984.20	(\$1,984.20)
[Blue]"\$"#,##0.00_);[Red]("\$"#,##0.00)	\$1,984.20	(\$1,984.20)
0%	198420%	-198420%
0.00%	4466984.03%	
0.00E+00	4.47E+04	
mm-dd-yy	4/18/22	
d-mmm-yy	18-Apr-22	
d-mmm	18-Apr	
mmm-yy	Apr-22	
h:mm AM/PM	8:10 PM	
h:mm:ss AM/PM	8:10:00 PM	
h:mm	20:10	
h:mm:ss	20:10:00	
m/d/yy h:mm	4/18/22 20:10	

The above is created with this code.

```

1 import openpyxl
2
3 formats = ['General',
4           '0',
5           '0.00',
6           '#,##0',
7           '#,##0.00',
8           '+ "$"#,##0_);("$"#,##0)',
9           '"$"#,##0_);[Red]("$"#,##0)',
10          '[Green]"$"#,##0_);[Red]("$"#,##0)',
11          '"$"#,##0.00_);("$"#,##0.00)',
12          '"$"#,##0.00_);[Red]("$"#,##0.00)',
13          '[Blue]"$"#,##0.00_);[Red]("$"#,##0.00)',
14          '0%',
15          '0.00%',
16          '0.00E+00',
17          'mm-dd-yy',

```

⁴Access the [spreadsheet](#) here.

```

18     'd-mmm-yy',
19     'd-mmm',
20     'mmm-yy',
21     'h:mm AM/PM',
22     'h:mm:ss AM/PM',
23     'h:mm',
24     'h:mm:ss',
25     'm/d/yy h:mm']
26
27 wb = openpyxl.Workbook()
28 ws = wb.active
29
30 ws['A1'] = 'Format String'
31 ws['B1'] = 'Example'
32 ws['C1'] = 'Example (negative)'
33 for key, fmt in enumerate(formats):
34
35     coord1 = 'A{}'.format(key + 2)
36     coord2 = 'B{}'.format(key + 2)
37     coord3 = 'C{}'.format(key + 2)
38
39     # plain format string
40     ws[coord1].value = fmt
41
42     # set a date or numeric value and apply format
43     if key > 11:
44         ws[coord2].value = dt.datetime(2022,4,18,20,10)
45     else:
46         ws[coord2].value = 1984.1984
47         ws[coord3].value = -1984.1984
48         ws[coord3].number_format = fmt
49     ws[coord2].number_format = fmt
50
51 wb.save("openpyxl_formats.xlsx")

```

This leaves [fonts](#), [borders](#), [alignment](#), and [fills](#) to be explored.

Here is an example with the Excel result further below.

```

1 import openpyxl
2 import openpyxl.styles as osty
3 wb = openpyxl.Workbook()
4 ws = wb.active
5
6 words = 'My cup runneth over.'.split(" ")
7 letters = 'ABCD'
8 for l, s in zip(letters, words):
9     ws['{}1'.format(l)].value = s
10
11 ws['A1'].font = osty.Font('Times New Roman', bold = True, color = '006400')
12
13 ws['B1'].alignment = osty.alignment.Alignment(horizontal = 'right',
14                                             vertical = 'top',
15                                             textRotation = 30)
16
17 ws['C1'].fill = osty.fills.PatternFill(fgColor = 'FFFFE3',
18                                       fill_type = 'solid')
19
20 medium = osty.borders.Side(border_style = 'medium',
21                           color = 'FF0000')
22 dotted = osty.borders.Side(border_style = 'dotted',
23                           color = '87E0FF')
24 ws['D1'].border = osty.borders.Border(top = medium,
25                                       bottom = medium,
26                                       left = dotted,
27                                       right = dotted)
28
29 wb.save("openpyxl_cup.xlsx")

```

My	cup	runneth	over.
----	-----	---------	-------

The above uses `PatternFill`. There's also `GradientFill`, used below. I'm not sure anyone has a use for it. Nonetheless:

```
1 wb = openpyxl.Workbook()
2 ws = wb.active
3
4 strings = ['Every man is important if he loses his life;',
5           'and every many is funny if he loses his hat and has to run after it.']
6
7 for key, s in enumerate(strings):
8     ws['A{}'.format(key+1)].value = s
9
10 for cell in ws['A']:
11     cell.font = osty.Font('Times New Roman', bold = True)
12
13 ws['A1'].fill = osty.fills.GradientFill(stop = ['1D97C1', 'FFFFFF'])
14 ws['A2'].fill = osty.fills.GradientFill(stop = ['FFFFFF', '008080'], degree = 90)
15
16 wb.save("openpyxl_chesterton.xlsx")
```

**Every man is important if he loses his life;
and every many is funny if he loses his hat and has to run after it.**

Conditional Formatting and Filtering

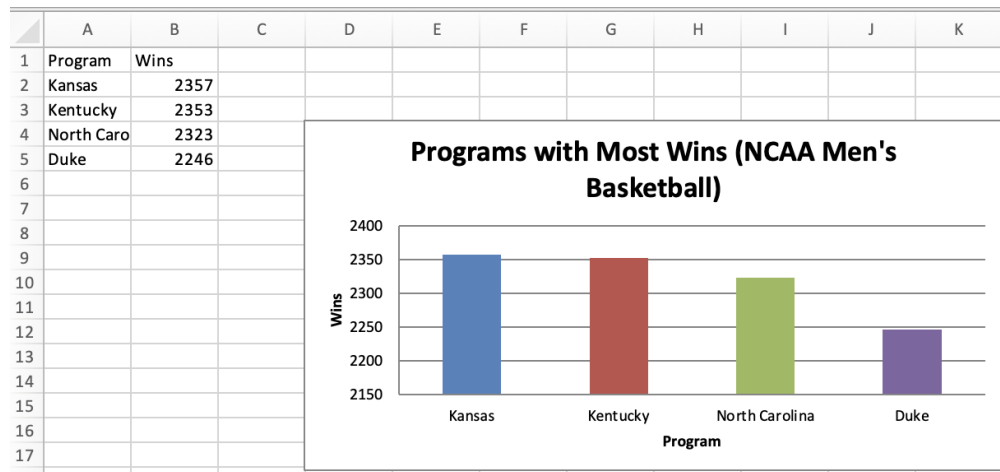
Next, we consider conditional formatting. In my experience, the vast majority of any Excel formatting I did was applying a filter and conditional formatting. Below, we take a `DataFrame` with cosine similarities and send it to a formatted Excel file.

```
1 from openpyxl.formatting.rule import ColorScaleRule
2 import numpy as np
3 np.random.seed(1)
4
5 # Data Generation
6 a = np.random.rand(n, n)
7 a = np.triu(a, k = 1) + np.triu(a).T
8 similarities = a - np.diag(np.diag(a)) + np.diag(np.ones(n))
9
10 # Send to DataFrame with index and headers
11 names = ['product{:0>2}'.format(x) for x in range(n)]
12 df = pd.DataFrame(similarities, index = names, columns = names)
13
14 # Create book and sheet
15 book = openpyxl.Workbook()
16 ws = book.active
17 ws.title = 'Product Similarities'
18
19 # Send data to worksheet
20 for r in dataframe_to_rows(df, index=True, header=True):
21     if r != [None]:
22         ws.append(r)
23
24 # Create and apply color rule
25 rule = ColorScaleRule(start_type='percentile', start_value=10, start_color='ea9999', # red
26                       mid_type='percentile', mid_value=50, mid_color='FFFFFF', # white
27                       end_type='percentile', end_value=90, end_color='b6d7a8') # green
28 ws.conditional_formatting.add('A1:AP42', rule)
29
30 # Add Filter over entire spreadsheet
31 ws.auto_filter.ref = ws.calculate_dimension()
```

[illegible]

Next, we'll add a bar chart using OpenPyXL. See [Zumstein 2021](#) Chapter 9 for adding matplotlib images to an Excel file. It's better to include an actual Excel chart as the Excel file user can make further modifications.

We get output like this.



Chapter 17

Applications

17.1 Wild Data Redux

In an earlier lecture, we worked with wild data for the first time by using the public Peloton API to look at the JSON file associated with a single workout. Now, we'll revisit that and go a bit further using our expanded toolkit.

Here are some new topics we'll learn about along the way.

1. Missing data and interpolation
2. Binning numerical data
3. Categorical data type

First, we'll replicate the output graph shown below in the Peloton workout details.



Anyone can access data from a public Peloton profile, provided they have a login. Still, we'll bypass that step in these notes and suppose we already have the JSON files.

```
1 with open('redux_workout.json') as json_file:
2     packets = json.load(json_file)
```

Check the keys with `.keys()` and you'll see we're interested in `packets['metrics']`. Further, this is a list, and we find the details for output at index level zero.

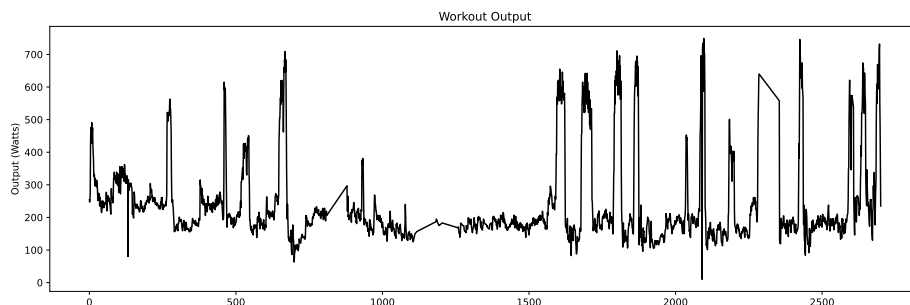
```
1 plt.subplots(figsize = (16,5))
2 plt.plot(pack2['metrics'][0]['values'])
3 plt.title("Workout Output")
4 plt.ylabel("Output (Watts)")
5 plt.savefig("redux_output1.pdf")
```



We should see that this does *not* match the workout details. For example, there is a straight line seen in the official details beginning at the 38 minute mark that does not show up in our plot. That's because there is missing data. We can verify this by checking the length of the packet array is not equal to the workout duration.

With more inspection, you will find `packets['seconds_since_pedaling_start']`, which is itself a list of seconds. These can be interpreted as the corresponding workout times for each of the metric values (this is more Peloton knowledge or educated guessing than something immediately obvious from the data). These are the seconds for which the metrics were successfully captured.

```
1 plt.subplots(figsize = (16,5))
2 plt.plot(packets['seconds_since_pedaling_start'], packets['metrics'][0]['values'],
3         color = 'black')
4 plt.title("Workout Output")
5 plt.ylabel("Output (Watts)")
6 plt.savefig("redux_output2.pdf")
```



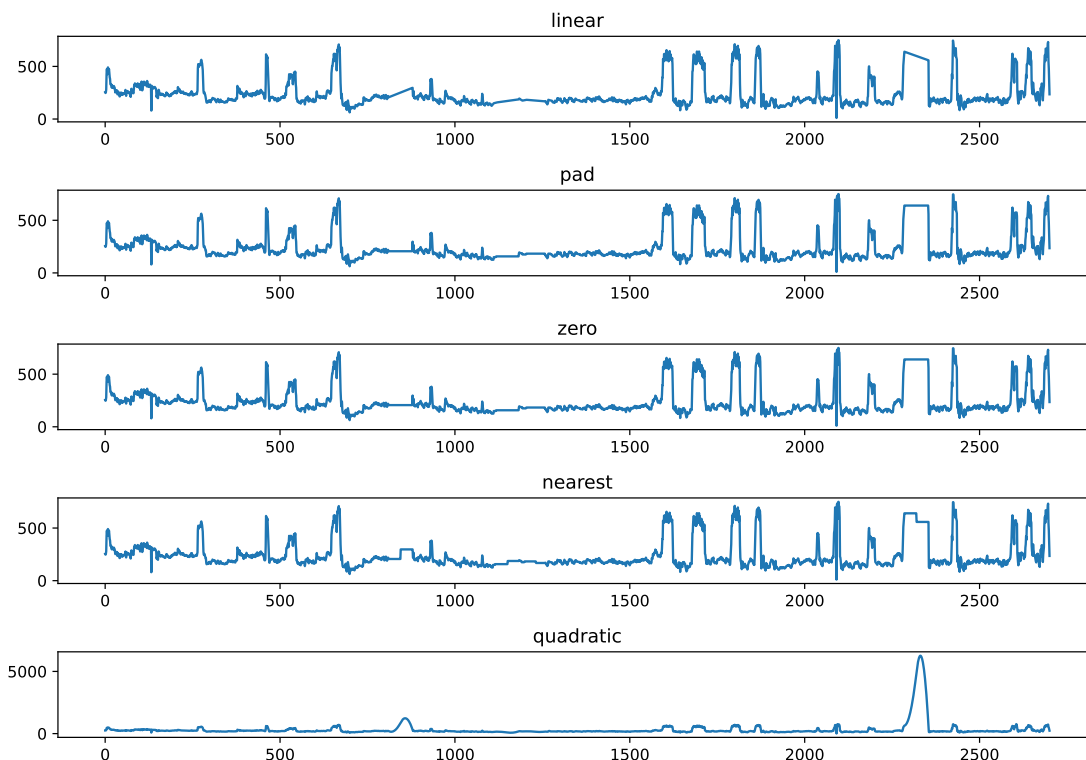
Now, this matches the Peloton workout details (even if it's not as pretty).

17.1.1 Missing Data

Above, we essentially interpolated the data linearly in the plot. That's because we drew a line between each value. So interpolation is done in a graph basically for free. In your data, you'll have to do more. The

simplest way is to use the Series method `interpolate()`. This takes a Series with missing values and fills in those missing values

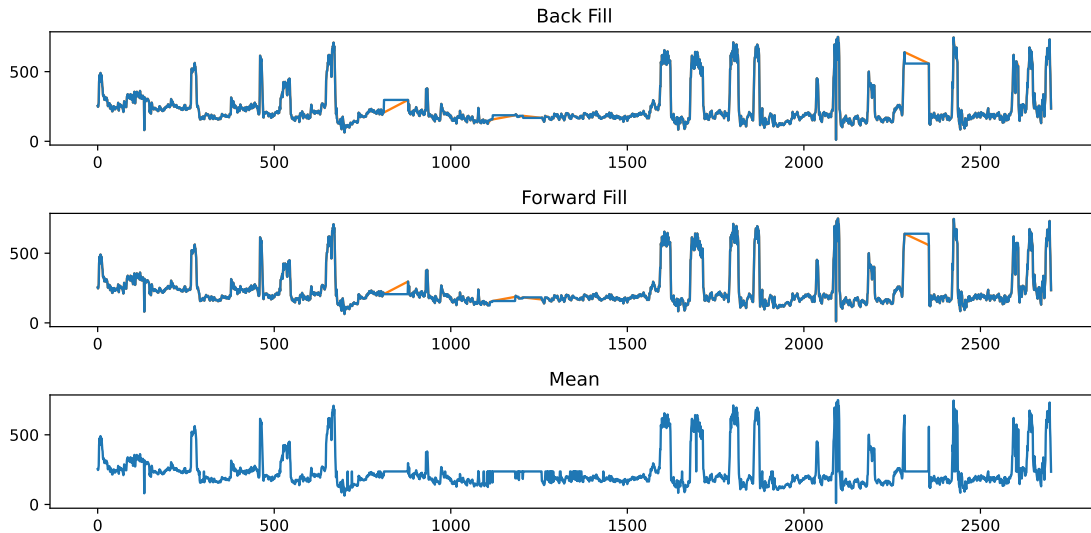
```
1 interp_methods = ['linear', 'pad', 'zero', 'nearest', 'quadratic']
2
3 fig, ax = plt.subplots(len(interp_methods), 1, figsize = (10,7))
4 for key, s in enumerate(interp_methods):
5     df['output'].interpolate(s).plot(ax = ax[key])
6     ax[key].set_title(s)
7 plt.tight_layout()
8 plt.savefig("sample_interpolations.pdf")
```



Notice the very different *y*-axis limits for the quadratic interpolation.

There are additional capabilities to use splines and more advanced ways of filling in the missing data. We won't cover that. Another low tech route is to use the `fillna()` Series method and related methods like `bfill()` and `ffill()`.

```
1 fig, ax = plt.subplots(3, 1, figsize = (10,5))
2
3 # I include linear interpolation in orange/C1 for reference
4 df['output'].interpolate().plot(ax = ax[0], color = 'C1')
5 df['output'].bfill().plot(ax = ax[0], color = 'C0')
6 ax[0].set_title("Back Fill")
7
8 df['output'].interpolate().plot(ax = ax[1], color = 'C1')
9 df['output'].ffill().plot(ax = ax[1], color = 'C0')
10 ax[1].set_title("Forward Fill")
11
12
13 mean_output = df.output.mean()
14 df['output'].fillna(mean_output).plot(ax = ax[2])
15 ax[2].set_title("Mean")
16
17 plt.tight_layout()
18 plt.savefig("fillna_examples.pdf")
```



17.1.2 Categorical

Next, we might prefer to look at lower resolution data. Cyclists sometimes don't look at the exact output value, but instead focus on staying in some range of values. This lends itself to [power zone training](#).

This means we might use (ordered) *categorical* data. This relies on the pandas `Categorical` type. Categories that are numerical intervals can be created with `cut()` and `qcut()`. You might also work with categories that are simply string names. `cut()` takes a sequence of numbers and places it in a numerical category. Below, we set the bins manually based off `zone_mins`. Then we can access the created categories with the `categories` attribute. There is also the accessor, `cat`, which works like the `str` string accessor. We use that when creating `zone_data`, transforming raw output to `zone/categorical` data.

```
1 ftp = 210
2 zone_mins = 0, .56*ftp, .76*ftp, .91*ftp, 1.06*ftp, 1.21*ftp, 1.51*ftp, 10**10
3
4 # Hacky way to get interval categories
5 cats = pd.cut([1], bins = zone_mins).categories
6
7 zone_data = df.output.astype('category').cat.set_categories(cats1)
8 zone_data.value_counts()
```

The categorical data type has the advantage of better performance and something like `value_counts()` behaves slightly different.

```
1 # Shows 0 frequencies for high zones
2 pd.cut([1], zone_mins).value_counts()
3
4 # Doesn't show high zones
5 pd.Categorical([pd.Interval(zone_mins[0], zone_mins[1])]).value_counts()
6
7 # Non-Peloton/non-interval example
8
9 a1 = pd.Series(['a'])
10 print(a1.value_counts())
11
12 a2 = pd.Series(['a']).astype('category')
13 a2 = a2.cat.set_categories(['a', 'b'])
14 print(a2.value_counts())
```

Dummy variables can be added to our dataset next. `get_dummies()` works with a array-like data, `Series`, and `DataFrames`. You might also call this one-hot encoding.

```
1 pd.get_dummies(zone_data)
```

17.2 Inference and Experiments

In this section, we put our skills together to analyze an A/B test. We'll use pandas, NumPy, and a couple new things.

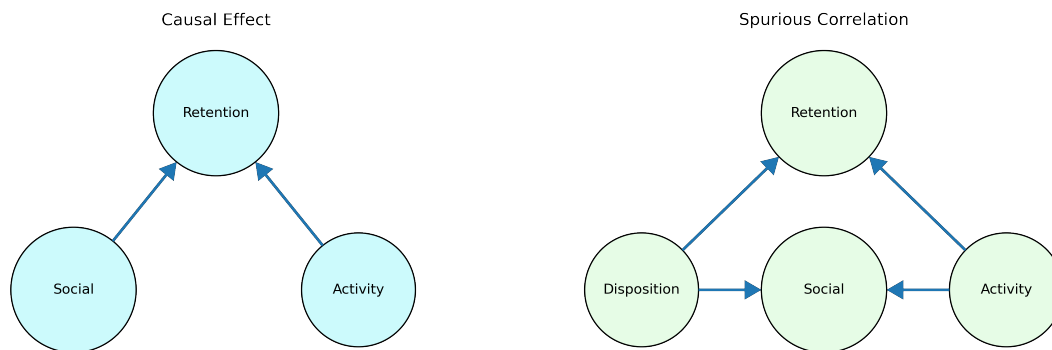
1. Use `statsmodels` for ordinary least squares.¹
2. Use `scipy` for a two-sample *t*-test.
3. Use `numpy` for randomization inference/permutation tests.

First, here's some motivation for why A/B tests are important and to create more student demand for the causal inference course I want to teach.

17.2.1 Motivation/Soapbox

Any economic or behavioral model proposes to describe the Data Generating Process (DGP), or whatever is really driving the outcomes we see in our non-experimental data. That is, we want to work out if X causes Y or X and Y have a common cause and the association between them is spurious. What would you say if you (plausibly) analyzed data that showed people wearing winter hats are colder than people not wearing winter hats? This is all to say that questions of causality are fraught with issues of bias. Hence the oft-repeated directive [not to confuse correlation with causality](#).

The diagrams below depict two different causal models. Suppose you work for a business with a subscription model. It might be Duolingo, Peloton, Netflix, etc. There is normal activity, which is simply using the product and there is social engagement, which might be liking or sharing posts on social media. A product manager might want to know the business impact of social engagement, and both models might create the same correlations between social and retention. A great analyst is sensitive to the idea that correlations are not proof of a causal relationship.



On the left, social engagement has a direct effect on the retention outcome. On the right, social engagement has no effect on the retention outcome. However, as we verify in the simulation below, both can produce the same regression output. We'll use OLS for simplicity.

We regress retention on social engagement and activity, first with data generated by the causal effect process (y_1), then the spurious correlation process (y_2).

```
1 import statsmodels.api as sm
2 import pandas as pd
3
4 n = 1000
5 noise = np.random.normal(0, .1, size = n)
6
```

¹You might notice VanderPlas 2016a uses `sklearn` for linear regression. Both are good. `sklearn` is more ML oriented and `statsmodels` is more stats oriented. As an example of this difference, `sklearn` will apply regularization to a logistic regression by default and `statsmodels` won't.

```

7 # for causal model
8 activity = np.random.normal(size = n)
9 social = np.random.normal(size = n)
10 intercept = np.ones(n)
11 disposition = np.random.normal(size = n)
12 social_alt = disposition + activity
13
14
15 df = pd.DataFrame()
16 for thing in ['activity', 'social', 'social_alt', 'disposition', 'intercept']:
17     df[thing] = globals()[thing]
18
19
20 # Causal Model
21 y1 = 1*social + 1*activity + noise
22 X1 = df[['activity', 'social', 'intercept']]
23 causal_model = sm.OLS(y1, X1).fit()
24 print(causal_model.summary())
25
26 # Spurious Model
27 y2 = 2*activity + 1*disposition + noise
28 X2 = df[['activity', 'social_alt', 'intercept']]
29 spurious_model = sm.OLS(y2, X2).fit()
30 print(spurious_model.summary())

```

The regression output is produced below. They're the same, and that's bad. That means our simple OLS model isn't causal evidence, because the same results could be observed with a non-causal data generating process.

Causal effect model:

Dep. Variable:	y	R-squared:	0.995
Model:	OLS	Adj. R-squared:	0.995
Method:	Least Squares	F-statistic:	9.894e+04
Date:	Mon, 18 Apr 2022	Prob (F-statistic):	0.00
Time:	14:11:29	Log-Likelihood:	878.45
No. Observations:	1000	AIC:	-1751.
Df Residuals:	997	BIC:	-1736.
Df Model:	2		

	coef	std err	t	P> t	[0.025	0.975]
activity	1.0034	0.003	313.158	0.000	0.997	1.010
social	1.0008	0.003	314.360	0.000	0.995	1.007
intercept	0.0024	0.003	0.760	0.447	-0.004	0.009

Spurious correlation model:

Dep. Variable:	y	R-squared:	0.998
Model:	OLS	Adj. R-squared:	0.998
Method:	Least Squares	F-statistic:	2.477e+05
Date:	Mon, 18 Apr 2022	Prob (F-statistic):	0.00
Time:	14:09:52	Log-Likelihood:	878.62
No. Observations:	1000	AIC:	-1751.
Df Residuals:	997	BIC:	-1737.
Df Model:	2		

	coef	std err	t	P> t	[0.025	0.975]
activity	1.0013	0.005	220.729	0.000	0.992	1.010
social_alt	1.0020	0.003	315.603	0.000	0.996	1.008
intercept	0.0026	0.003	0.811	0.418	-0.004	0.009

This is frustrating, but it's why experiments hold a special place in economics and data science. For more on experiments and inference, you might check out Luca and Bazerman 2020 for commentary on experiments and Cunningham 2021 for a text on causal inference with Python code samples.

17.2.2 Experiments

t-tests

A simple A/B test is almost always analyzed by conducting a *t*-test. And typically, we have a two-sample test. For this, we can use [SciPy](#) and its stats submodule. We'll use the `scipy.stats.ttest_ind()` function to run a Welch's *t*-test.

```
1 import scipy.stats
2 import numpy as np
3 np.random.seed(1)
4
5 n = 1000 # observations
6 n_sims = 5000 # simulations
7 n_significant_results = 0
8 for _ in range(n_sims):
9
10     # exactly the same!
11     g1 = np.random.binomial(n = 1, p = 0.5, size = n)
12     g2 = np.random.binomial(n = 1, p = 0.5, size = n)
13
14     n_significant_results += 1 * (scipy.stats.ttest_ind(g1,g2).pvalue <= 0.05)
15
16
17 print(n_significant_results/n_sims)
```

We get false positives 5.38% of the time.

Lady Tasting Tea

Now, we'll move toward calculating exact *p*-values using simulations. To start, we consider the classic [lady tasting tea](#) experiment.

Irving Fisher calculated that a random guesser could get the correct result only one time out of 70. Below, we verify this with a simulation.

```
1 import numpy as np
2 np.random.seed(24)
3
4 # without loss of generality, let truth be...
5 truth = np.array([1,1,1,1, 0,0,0,0])
6
7
8 n = 100_000
9 p_value = len([x for x in range(n) if np.random.permutation(truth)[0:4].sum() == 4]) / n
10
11 print(p_value)
```

Randomization Inference

Randomization inference goes just a bit beyond what's done in the lady tasting tea example. We have a single observed treatment effect and then we simulate many more counterfactual treatment effects by permuting the treatment and control assignments. The permutations simulate the kinds of treatment effects we'd observe under the null hypotheses. If we calculate how often those treatment effects are more extreme than what we actually observed, we're left with an *exact p*-value.

```
1 users = pd.read_csv('users_l12.csv')
2 engagement = pd.read_csv('engagement_l12.csv')
3
4 # Simulate 1000 alternate labels
5 for i in range(1000):
6     users['alt{}'.format(i)] = np.random.permutation(users.assignment)
7
8 df2 = engagement.merge(users, on = 'user_id')
9
10 # homework to improve this step
```

```
11 treatment_effects = list()
12 for column in [x for x in list(df2) if 'alt' in str(x)]:
13     sums = df2.groupby(column).minutes_engaged.sum()
14     te = sums['treatment'] - sums['control']
15     treatment_effects.append(te)
16
17
18 # Make Plot
19 plt.hist(treatment_effects, bins = 30)
20
21 # actual treatment effect
22 grouped = df2.groupby('assignment').minutes_engaged.sum()
23 te_true = grouped['treatment'] - grouped['control']
24 plt.axvline(te_true, color = 'black')
25
26 plt.show()
27
28 abs_values = np.abs(treatment_effects)
29
30 p_value = len(abs_values[abs_values > treatment_effect()]) / 1000
31
32 print(p_value)
33
34
35 # compare to t-test
36
37 from scipy.stats import ttest_ind
38
39 treatment_values = df[df.assignment == 'treatment'].minutes_engaged.values
40 control_values = df[df.assignment == 'control'].minutes_engaged.values
41
42 test = ttest_ind(treatment_values, control_values, equal_var = False)
43 print('t-test p value', test.pvalue)
```

Bibliography

- Antao, T. R. (2022). *Fast python for data science*. Manning Publications.
- Clark, A. (2022). *Matplotlib for storytellers*.
- Cunningham, S. (2021). *Causal inference: The mixtape*. Yale University Press.
- Gaddis, T. (2018). *Starting out with python*. Pearson.
- Long, T. (2021). *Good code, bad code*. Manning Publications.
- Lubanovic, B. (2019). *Introducing python: Modern computing in simple packages*. O'Reilly Media.
- Luca, M., & Bazerman, M. H. (2020). Want to make better decisions? start experimenting. *MIT Sloan Management Review*, 61(4), 67–73.
- McKinney, W. (2017). *Python for data analysis: Data wrangling with pandas, numpy, and ipython*. O'Reilly Media, Inc.
- Ramalho, L. (2015). *Fluent python: Clear, concise, and effective programming*. O'Reilly Media.
- Reitz, K., & Schlusser, T. (2016). *The hitchhiker's guide to python: Best practices for development*. O'Reilly Media.
- Rougier, N. (2021). *Scientific visualization: Python+ matplotlib*.
- VanderPlas, J. (2016a). *Python data science handbook: Essential tools for working with data*. O'Reilly Media.
- VanderPlas, J. (2016b). *A whirlwind tour of python*. O'Reilly Media.
- Zumstein, F. (2021). *Python for excel*. O'Reilly Media.