

YData: Introduction to Data Science



Class 13: loops and writing functions

Overview

Quick review of maps

For loops

Conditional statements

Writing functions

If there is time: text manipulation



Announcement: Homework 6

Homework 6 has been posted!

It is due on Gradescope on **Sunday March 3rd at 11pm**

Reminder: class project

The final project is a **6-10 page** Jupyter notebook report where you analyze your own data to address a question that you find interesting!

Start thinking about what questions you want to examine and start looking into getting data

- A few sources for data sets are listed on Canvas

A **polished** draft of the project is due on **April 7th**



Midterm exam

Thursday March 7th **in person** during regular class time

- Exam is on paper

As part of homework 6, you will post a practice problem to Ed

- Ideally do this soon
- I will take one of these problems and put it on the exam

A practice exam has been posted



Midterm exam “cheat sheet”

You are allowed an exam “cheat sheet”

One page, double sided, that contains **only code**

- No code comments allowed
- E.g., `sns.catplot(data = , x = , y = , hue = , kind = "strip"/"swarm")`

Cheat sheet must be on a regular 8.5 x 11 piece of paper

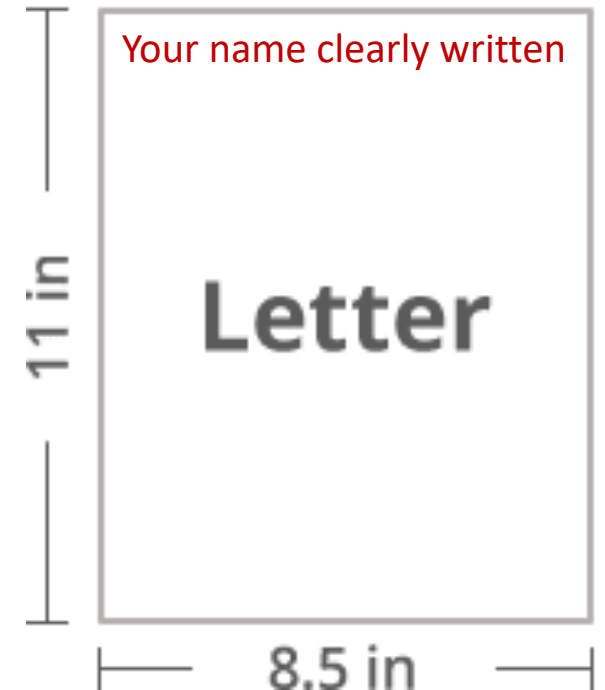
- Your name on the upper left of both sides of the paper

Strongly recommend making a typed list of all functions discussed in class and on the homework

- This will be useful beyond the exam

You must turn in your cheat sheet with the exam

- Failure to do so will result in a 20 point deduction



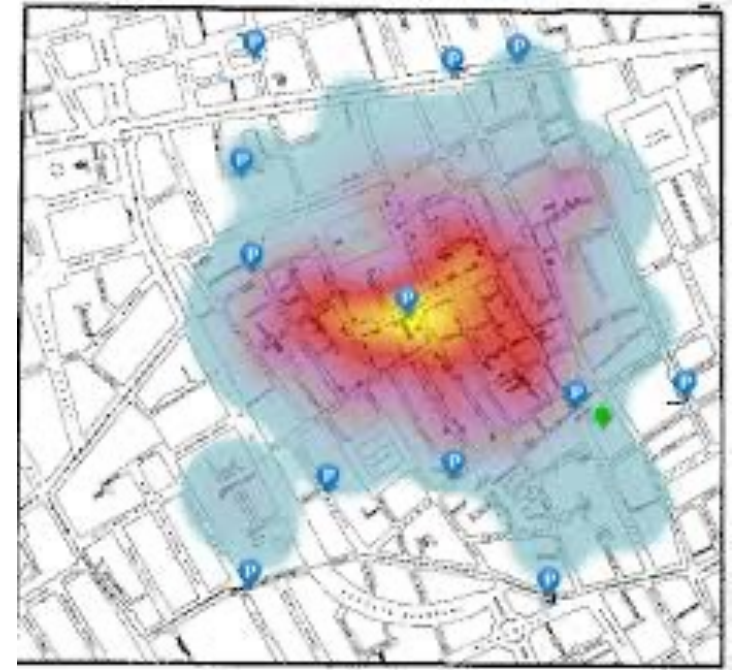
Quick review of mapping

Quick review of maps

Visualizing data on a map can be a powerful way to see spatial trends

We can create maps in Python using geopandas DataFrames

- Like regular DataFrames with an additional geometry column that has Shaply objects



John Snow's ghost map (1854)

	key_comb_drvr	geometry
0	M11551	POINT (117.525391 34.008926)
1	M17307	POINT (86.51248 30.474344)
2	M19584	POINT (89.537415 37.157627)

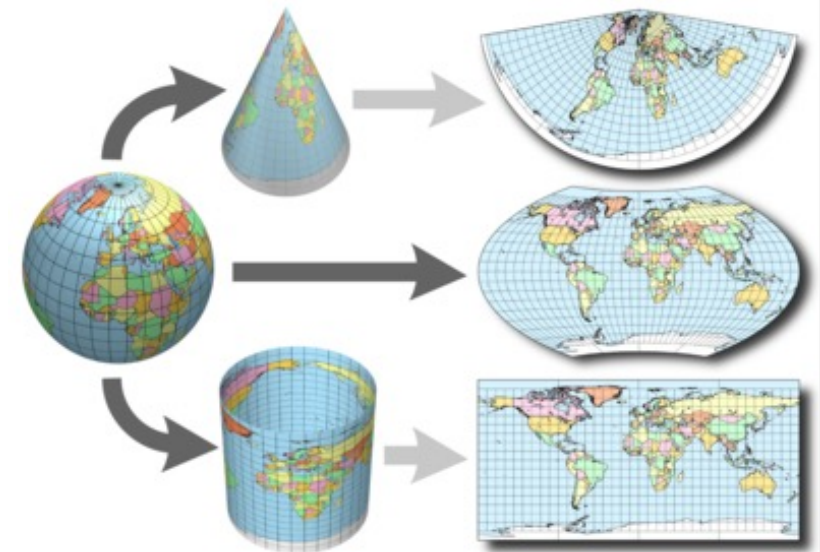
Review: CRSs and map projections

A coordinate reference system (CRS) is a framework used to precisely measure locations on the surface of the Earth as coordinates

- Needed for aligning different layers on maps

There are many map projections to display Earth's 3D structure on a 2D map surface.

- **Mercator projection** keeps angles intact
- **Eckert IV projection** keeps the size of land areas intact



WHAT YOUR FAVORITE MAP PROJECTION SAYS ABOUT YOU

MERCATOR



YOU'RE NOT REALLY INTO MAPS.

VAN DER GRINTEN



YOU'RE NOT A COMPLICATED PERSON. YOU LOVE THE MERCATOR PROJECTION; YOU JUST WISH IT WEREN'T SQUARE. THE EARTH'S NOT A SQUARE, IT'S A CIRCLE. YOU LIKE CIRCLES. TODAY IS GONNA BE A GOOD DAY!

HOB0-DYER



YOU WANT TO AVOID CULTURAL IMPERIALISM, BUT YOU'VE HEARD BAD THINGS ABOUT GAIL-PETERS. YOU'RE CONFLICT-AVERSE AND BUY ORGANIC. YOU USE A RECENTLY-INVENTED SET OF GENDER-NEUTRAL PRONOUNS AND THINK THAT WHAT THE WORLD NEEDS IS A REVOLUTION IN CONSCIOUSNESS.

PLATE CARRÉE (EQUIRECTANGULAR)



YOU THINK THIS ONE IS FINE. YOU LIKE HOW X AND Y MAP TO LATITUDE AND LONGITUDE. THE OTHER PROJECTIONS OVERCOMPLICATE THINGS. YOU WANT ME TO STOP ASKING ABOUT MAPS SO YOU CAN ENJOY DINNER.

ROBINSON



YOU HAVE A COMFORTABLE PAIR OF RUNNING SHOES THAT YOU WEAR EVERYWHERE. YOU LIKE COFFEE AND ENJOY THE BEATLES. YOU THINK THE ROBINSON IS THE BEST-LOOKING PROJECTION, HANDS DOWN.

DYMAXION



YOU LIKE ISAAC ASIMOV, XML, AND SHOES WITH TOES. YOU THINK THE SEAWAY GOT A BAD RAP. YOU OWN 3D GOGGLES, WHICH YOU USE TO VIEW ROTATING MODELS OF BETTER 3D GOGGLES. YOU TYPE IN DVDRMK.

A GLOBE!



YES, YOU'RE VERY CLEVER.

WATERMAN BUTTERFLY



REALLY? YOU KNOW THE WATERMAN? HAVE YOU SEEN THE 1909 CAHILL MAP ITS BASED— ... YOU HAVE A FRAMED REPRODUCTION AT HOME?! WHOA ... LISTEN, FORGET THESE QUESTIONS. ARE YOU DOING ANYTHING TONIGHT?

WINKEL-TRIPLE



NATIONAL GEOGRAPHIC ADOPTED THE WINKEL-TRIPLE IN 1998, BUT YOU'VE BEEN A WAT FAN SINCE LONG BEFORE 'NAT'GEO' SHOWED UP. YOU'RE WORRIED IT'S GETTING PLAYED OUT, AND ARE THINKING OF SWITCHING TO THE KAVRANSKY. YOU ONCE LEFT A PARTY IN DISGUST WHEN A GUEST SHOWED UP WEARING SHOES WITH TOES. YOUR FAVORITE MUSICAL GENRE IS "POST-".

GOODE HOMOLOGINE



THEY SAY MAPPING THE EARTH ON A 2D SURFACE IS LIKE FLATTENING AN ORANGE PEEL, WHICH SEEMS EASY ENOUGH TO YOU. YOU LIKE EASY SOLUTIONS. YOU THINK WE WOULDN'T HAVE SO MANY PROBLEMS IF WE'D JUST ELECT *ADORABLE* PEOPLE TO CONGRESS INSTEAD OF POLITICIANS. YOU THINK AIRLINES SHOULD JUST BUY ROOD FROM THE RESTAURANTS NEAR THE GATES AND SERVE *JAM* ON BOARD. YOU CHANGE YOUR OILS OIL, BUT SECRETLY WONDER IF YOU REALLY *NEED* TO.

PEIRCE QUINCUNCIAL



YOU THINK THAT WHEN WE LOOK AT A MAP, WHAT WE REALLY SEE IS OURSELVES. PETER YOU FIRST SAW *INCEPTION*, YOU SAT SILENT IN THE THEATER FOR SIX HOURS. IT BREAKS YOU OUT TO REALIZE THAT EVERYONE AROUND YOU HAS A SKELETON INSIDE THEM. YOU *HAVE* REALLY LOOKED AT YOUR HANDS.

GAIL-PETERS



I HATE YOU.

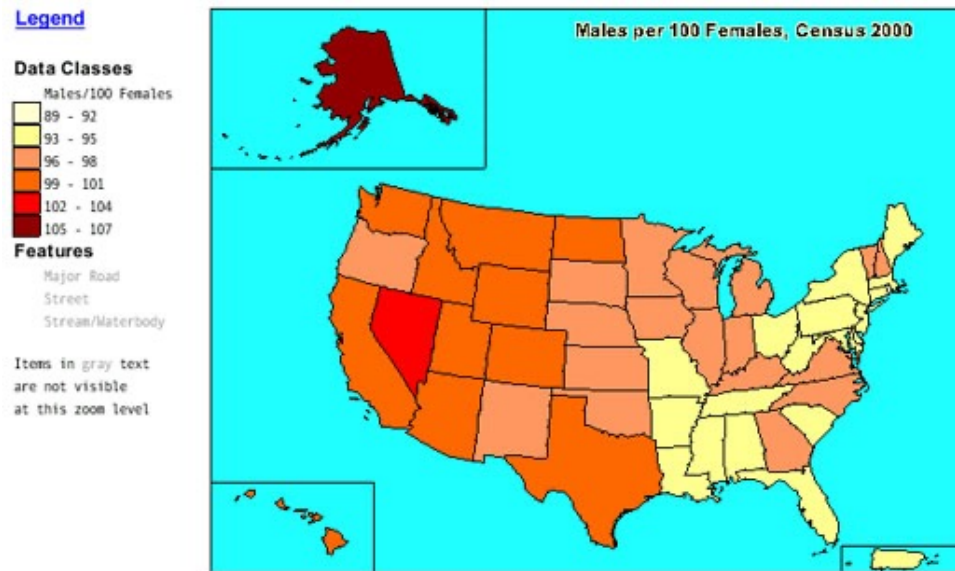
Review: Choropleth and Isopleth maps

Choropleth maps: shades/colors in predefined areas based on properties of a variable

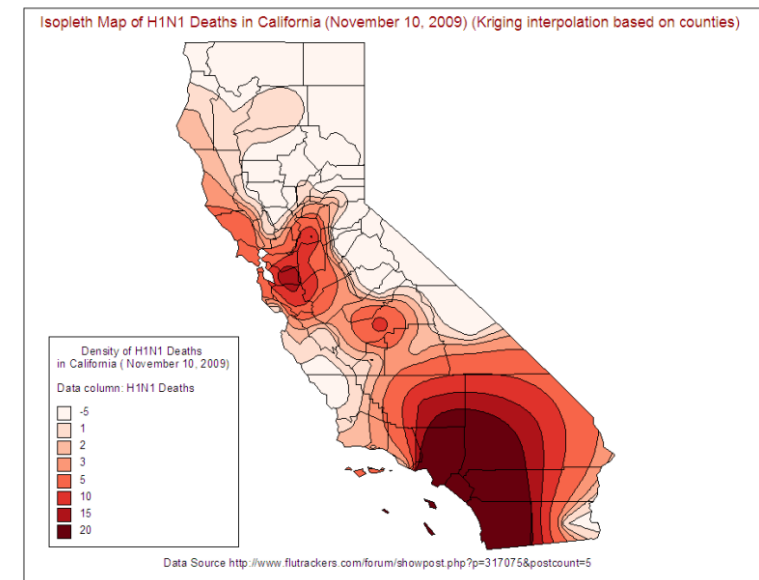
- We can then use the `gpd.plot(column =)` method to create choropleth maps

Isopleth maps: creates regions based on constant values

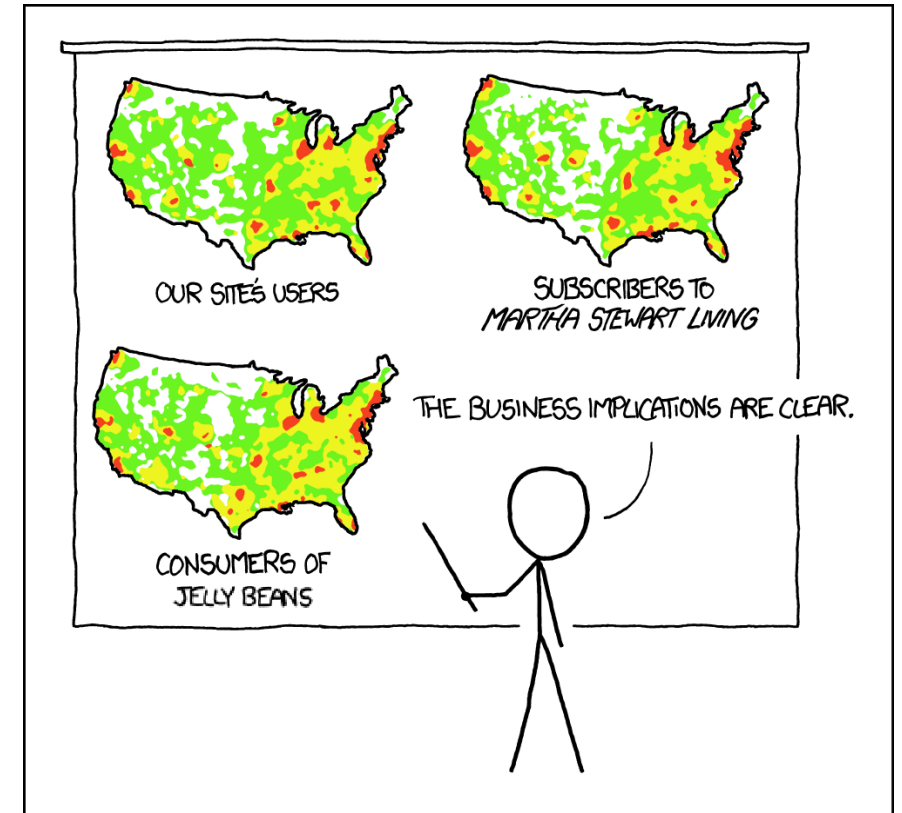
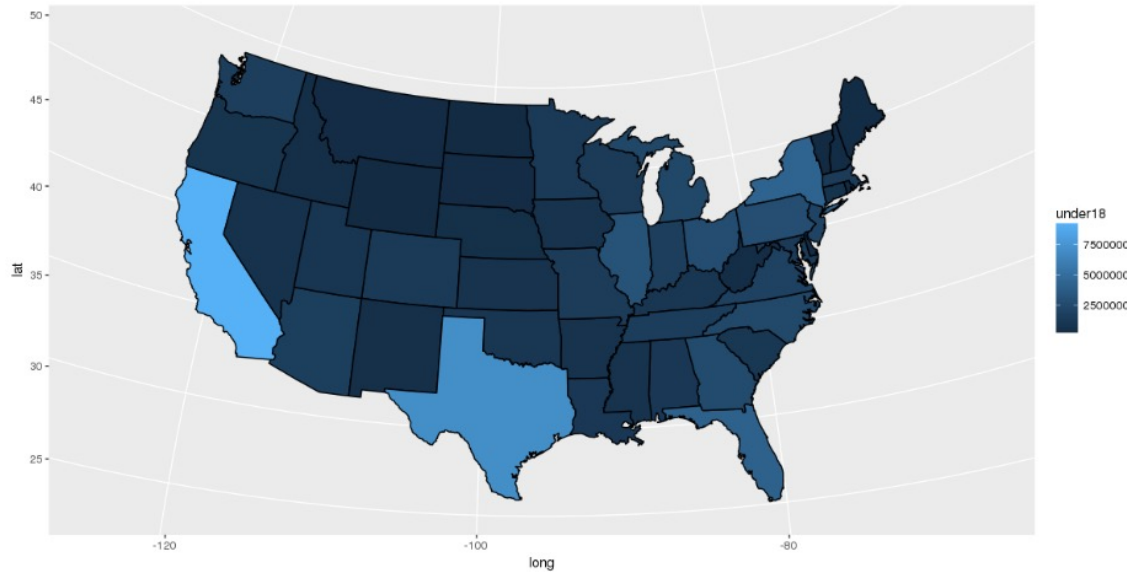
Choropleth map



Isopleth map

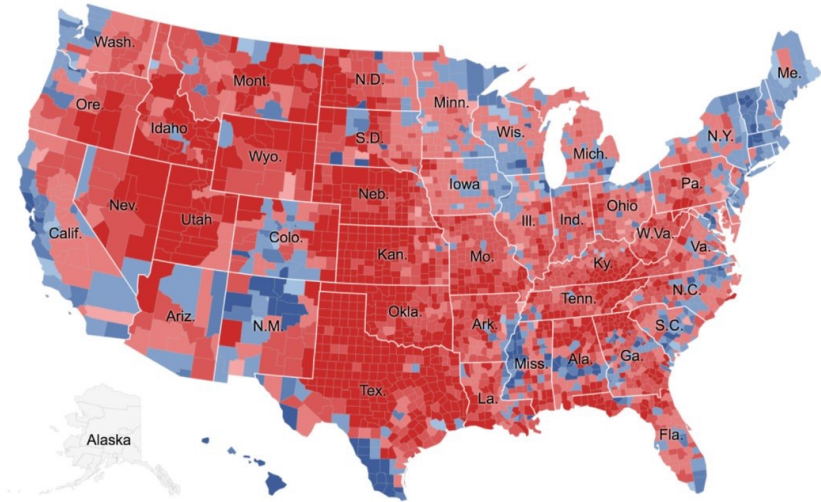


Review: Pet Peeve #208

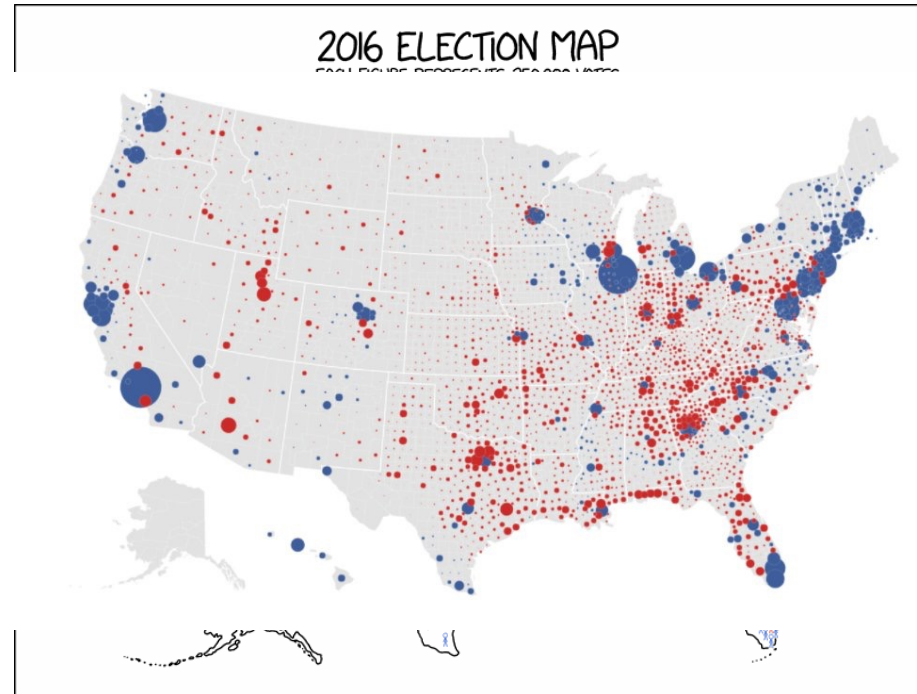


PET PEEVE #208:
GEOGRAPHIC PROFILE MAPS WHICH ARE
BASICALLY JUST POPULATION MAPS

Choropleth maps can be misleading



Looks like most of the country
voted republican



Let's look at a brief demo in Jupyter!

Loops

For loops

For loops repeat a process many times, iterating over a sequence of items

- Often we are iterating over an array of sequential numbers

```
animals = ["cat", "dog", "bat"]
```

```
for creature in animals:
```

```
    print(creature)
```

```
for i in np.arange(4):
```

```
    print(i**2)
```


Review: ranges

A range gives us a sequence of consecutive numbers

An sequence of increasing integers from 0 up to *end* - 1

- `range(end)`

An sequence of increasing integers from *start* up to *end* - 1

- `range(start, end)`

A sequence with step between consecutive values

- `range(start, end, step)`

The range always includes start but excludes end



Let's explore this in Jupyter!

Conditional statements

Review: comparisons

We can use mathematical operators to compare numbers and strings

- Results return Boolean values **True** and **False**

Comparison	Operator	True example	False Example
Less than	<	2 < 3	2 < 2
Greater than	>	3 > 2	3 > 3
Less than or equal	<=	2 <= 2	3 <= 2
Greater or equal	>=	3 >= 3	2 >= 3
Equal	==	3 == 3	3 == 2
Not equal	!=	3 != 2	2 != 2

We can also make comparisons across elements in an array

Conditional statements

Conditional statements control the sequence of computations that are performed in a program

We use the keyword **if** to begin a conditional statement to only execute lines of code if a particular condition is met.

We can use **elif** to test additional conditions

We can use an **else** statement to run code if none of the if or elif conditions have been met.

```
num = 5
if num == 1:
    print("Monday")
elif num == 2:
    print("Tuesday")
elif num == 3:
    print("Wednesday")
elif num == 4:
    print("Thursday")
elif num == 5:
    print("Friday")
elif num == 6:
    print("Saturday")
elif num == 7:
    print("Sunday")
else:
    print("Invalid input")
```

Let's explore this in Jupyter!

Defining functions

Writing functions

We have already used many functions that are built into Python or are imported from different modules/packages.

Examples...???

- `sum()`
- `statistics.mean()`
- `np.diff()`
- etc.

Let's now write our own functions!



Def statements

User-defined functions give names to blocks of code

```
def spread (values) :
```

Let's explore this in Jupyter!

text

MaNiPuLaTiOn

Text manipulation

80% of a Data Scientists time is cleaning data

- Text manipulation is a big part of cleaning data

20% of a Data Scientists time is complaining about cleaning data

Python has many string methods that are useful for manipulating text and cleaning data!

Text manipulation: capitalization

Some of the simplest string methods involve changing capitalization.

Changing capitalization can be useful when joining DataFrames

- i.e., if they key values are the same, but the values have different capitalization
 - For example, joining different countries, but in one DataFrame the country names are capitalized and in the other they are not



Text manipulation: capitalization

Python strings have a number of methods to change the capitalization of words including:

- `.capitalize()`: Converts the first character to upper case
- `.lower()`: Converts a string into lower case
- `.upper()`: Converts a string into upper case
- `.title()`: Converts the first character of each word to upper case
- `.swapcase()`: Swaps cases, lower case becomes upper case and vice versa

Let's explore this in Jupyter!

Text manipulation: string padding

Often we want to remove extra spaces (called "white space") from the front or end of a string.

Conversely, sometimes we want to add extra spaces to make a set of strings the same length

- This is known as "string padding"

Python strings have a number of methods that can pad/trim strings including:

- `strip()`: Returns a trimmed version of the string (i.e., with no leading or trailing white space).
 - Also, `rstrip()` and `lstrip()`: Returns a right/left trim version of the string
- `center(num)`: Returns a centered string (with equal padding on both sides)
 - Also `ljust(num)` and `rjust(num)`: Returns a right justified version of the string
- `zfill(num)`: Fills the string with a specified number of 0 values at the beginning

Let's explore this in Jupyter!

Text manipulation: checking string properties

There are also many functions to check properties of strings including:

- `isalnum()`: Returns True if all characters in the string are alphanumeric
- `isalpha()`: Returns True if all characters in the string are in the alphabet
- `isnumeric()`: Returns True if all characters in the string are numeric

- `isspace()`: Returns True if all characters in the string are whitespaces

- `islower()`: Returns True if all characters in the string are lower case
- `isupper()`: Returns True if all characters in the string are upper case
- `istitle()`: Returns True if the string follows the rules of a title

Let's explore this in Jupyter!

Text manipulation: splitting and joining strings

There are several methods that can help us join strings that are contained into a list into a single string, or conversely, parse a single string into a list of strings. These include:

- `split(separator_string)`: Splits the string at the specified separator, and returns a list
- `splitlines()`: Splits the string at line breaks and returns a list
- `join(a_list)`: Converts the elements of an iterable into a string

Let's explore this in Jupyter!

Text manipulation: finding and replacing substrings

Some methods for locating a substring within a larger string include:

- `count(substring)`: Returns the number of times a specified value occurs in a string
- `rfind(substring)`: Searches the string for a specified value and returns the last position of where it was found.
- `startswith(substring)`: Returns true if the string starts with the specified value
- `endswith(substring)` : Returns true if the string ends with the specified value
- `replace(original_str, replacement_str)`: Replace a substring with a different string.

Let's explore this in Jupyter!

Text manipulation: filling in strings with values

There are a number of ways to fill in strings parts of a string with particular values.

Perhaps the most useful is to use "f strings", which have the following syntax such as:

- `value_to_fill = "my_value"`
- `f"my string {value_to_fill} will be filled in"`

Let's explore this in Jupyter!

Regular expressions!



`/(reg)ex/`

Regular expressions

Regular expressions are string that allow you find more complex patterns in pieces of text

- They are powerful although can be a bit hard to read

`[^]*?@[^]*?\.[^]*`

To use regular expressions in Python we can import the re module

```
import re
```

We can check if a piece of text contains a particular substring by converting the output of `re.match()` method into a Boolean

```
bool(re.match("regular_expression", "piece_of_text"))
```

Regular expressions

[] means match anything in the range inside the braces

- "ch[io]mp" matches "chimp" and "chomp"

Note: if the ^ appears inside square braces it means **not**

- ^[^aeiou] matches words that don't start with a lower case vowel

The following are special regular expression characters that are reserved:

. * \ \$ { } [] ^ ?

Let's explore this in Jupyter!

Regular expressions

- (period) matches any single character
 - `bool(re.match("m.ss", "mess"))`
- * means match 0 or more of the preceding character
 - `bool(re.match("xy*z", "xz"))`
- + means match 1 or more of the preceding character
 - `bool(re.match("xy+z", "xz"))`

will the following match?

- `bool(re.match(".*a.*e", "pineapple"))`

Example

```
phone_strings = [ "apple",  
                  "219 733 8965",  
                  "329-293-8753",  
                  "Work: 579-499-7527",  
                  "Home: 203.867.9305"]
```

The phone number can be matched with the regular expression:

```
".*([2-9][0-9]{2})[- .]([0-9]{3})[- .]([0-9]{4})"
```

Let's explore this in Jupyter!

Escape sequences

In regular expressions a period (.) means any character

- So how can you detect if a period is in a string?

Escape sequences in R start with two slashes `\\` and cause the next character to be treated literally rather than as a special character

- To match a period we use `\\.`
- To match a \$ symbol we use `\\$`

Example

- `bool(re.match(".*\\$100", "Joanna has $100 and Chris has $0"))`

Character classes

Other special characters are also designated by using a double slash first

- `\s` space
- `\n` new line or also `\r`
- `\t` tab

Let's explore this in Jupyter!