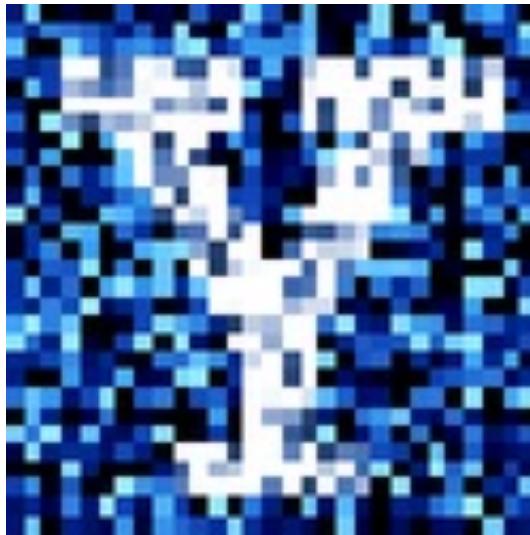


# YData: Introduction to Data Science



Class 13: Data visualization with seaborn

# Overview

Quick review of data visualization with  
matplotlib

Data visualization with seaborn

If there is time

- Text manipulation



# Announcement: Homework 6

Homework 6 has been posted!

It is due on Gradescope on **Sunday March 5<sup>th</sup>** at 11pm

- Be sure to mark each question on Gradescope along with the **page that has the answers!**

How was homework 5?

# Midterm exam

Thursday March 9<sup>th</sup> **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 will be posted soon



# Midterm exam “cheat sheet”

You are allowed an exam “cheat sheet”

One page, double sided, that contains only code

- No code comments allowed

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

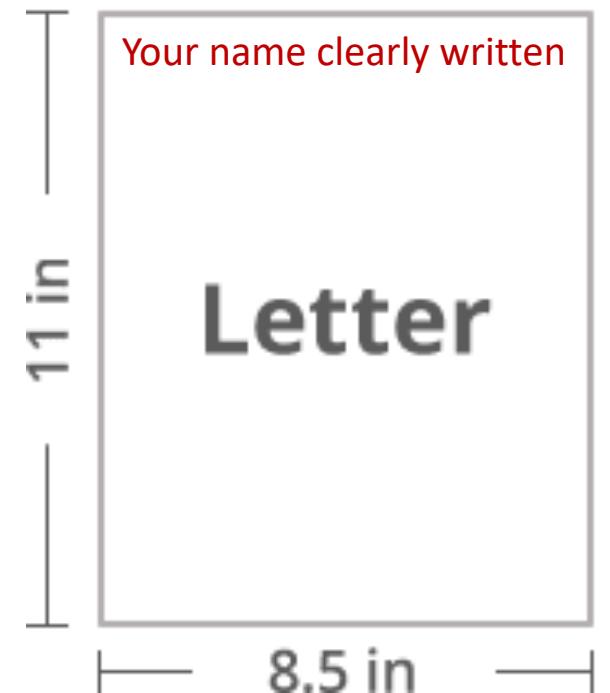
- Your name on the upper left of both sides of the paper

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



# Review of data visualization!



# Review of data visualization!

*Statistical projections which **speak to the senses without fatiguing the mind**, possess the advantage of fixing the attention on a great number of important facts.*

*—Alexander von Humboldt, 1811*

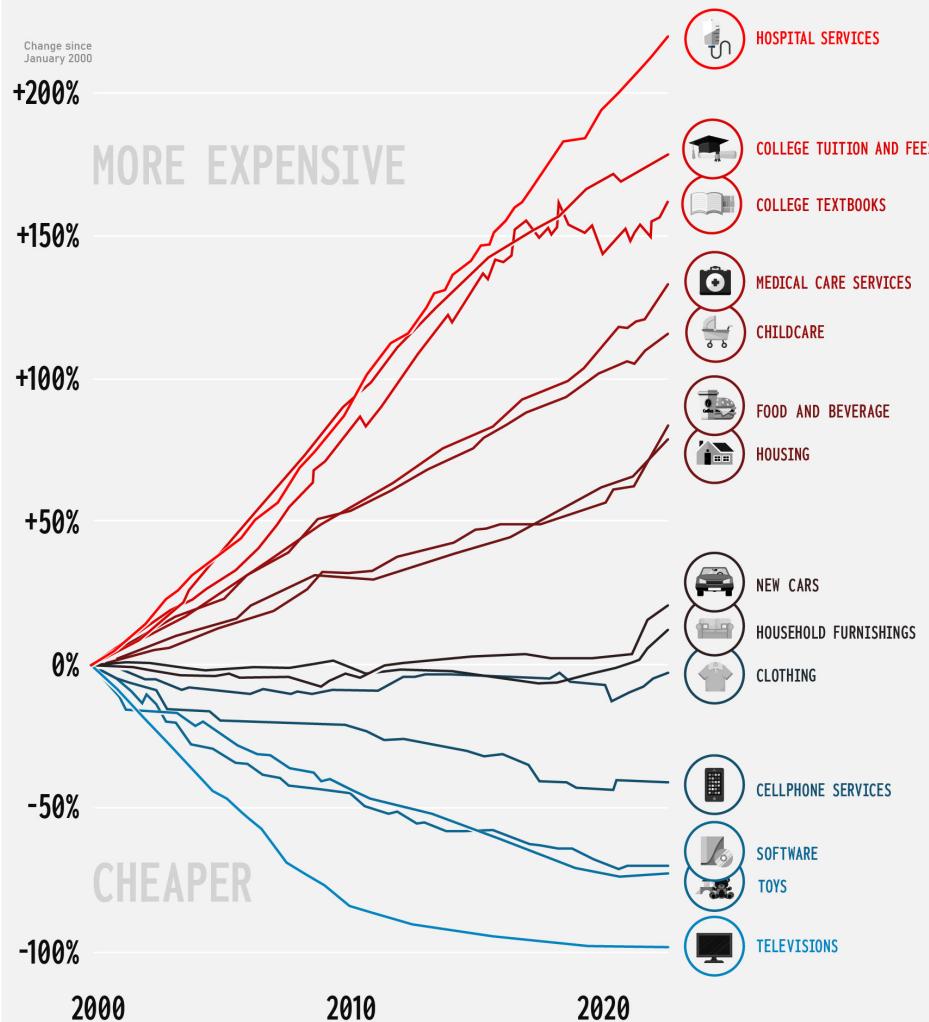


Let's take a few minutes to explore the data visualizations you found and that you created as part of homework 5...

# Line graphs

# Price Changes OF CONSUMER GOODS AND SERVICES

Broadly speaking, price levels have increased by 74% since 2000. That said, the actual numbers vary wildly depending on the type of good or service. Many consumer goods like **toys** and **TVs** have gotten cheaper, while critical categories like **healthcare** and **education** have skyrocketed.



Source: Bureau of Labor Statistics Original design and concept by Mark J. Perry, Senior Fellow Emeritus, AEI

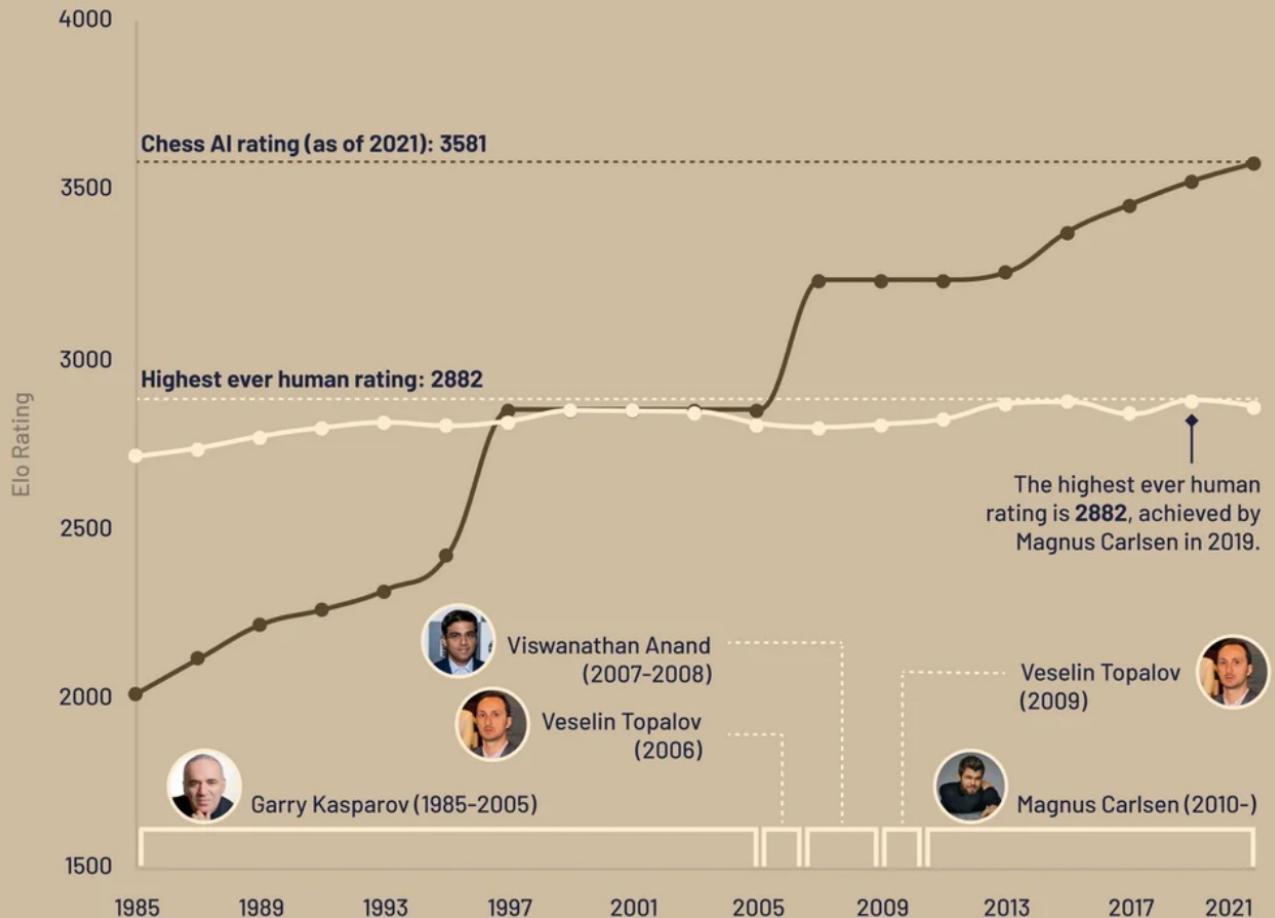
# Chess



## Artificial Intelligence vs Human

Created by  
genuine  
impact

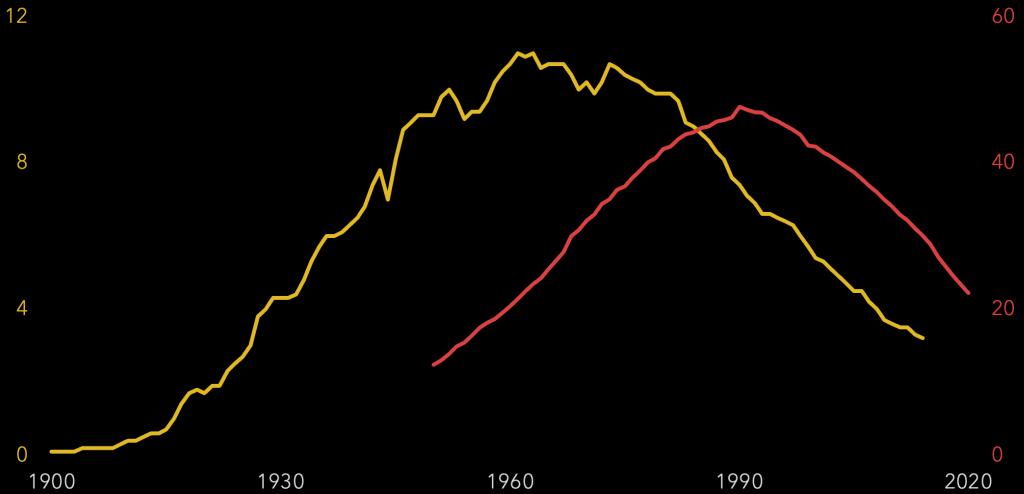
The Elo Rating system is used to measure chess ability, and is based on game results. A higher rating = greater ability.



Source: Our World in Data, FIDE

More charts: [genuineimpact.substack.com](https://genuineimpact.substack.com)

USA:  
Cigarettes Sold (per adult\* per day)  
Lung Cancer Death Rate (per 100,000 people)



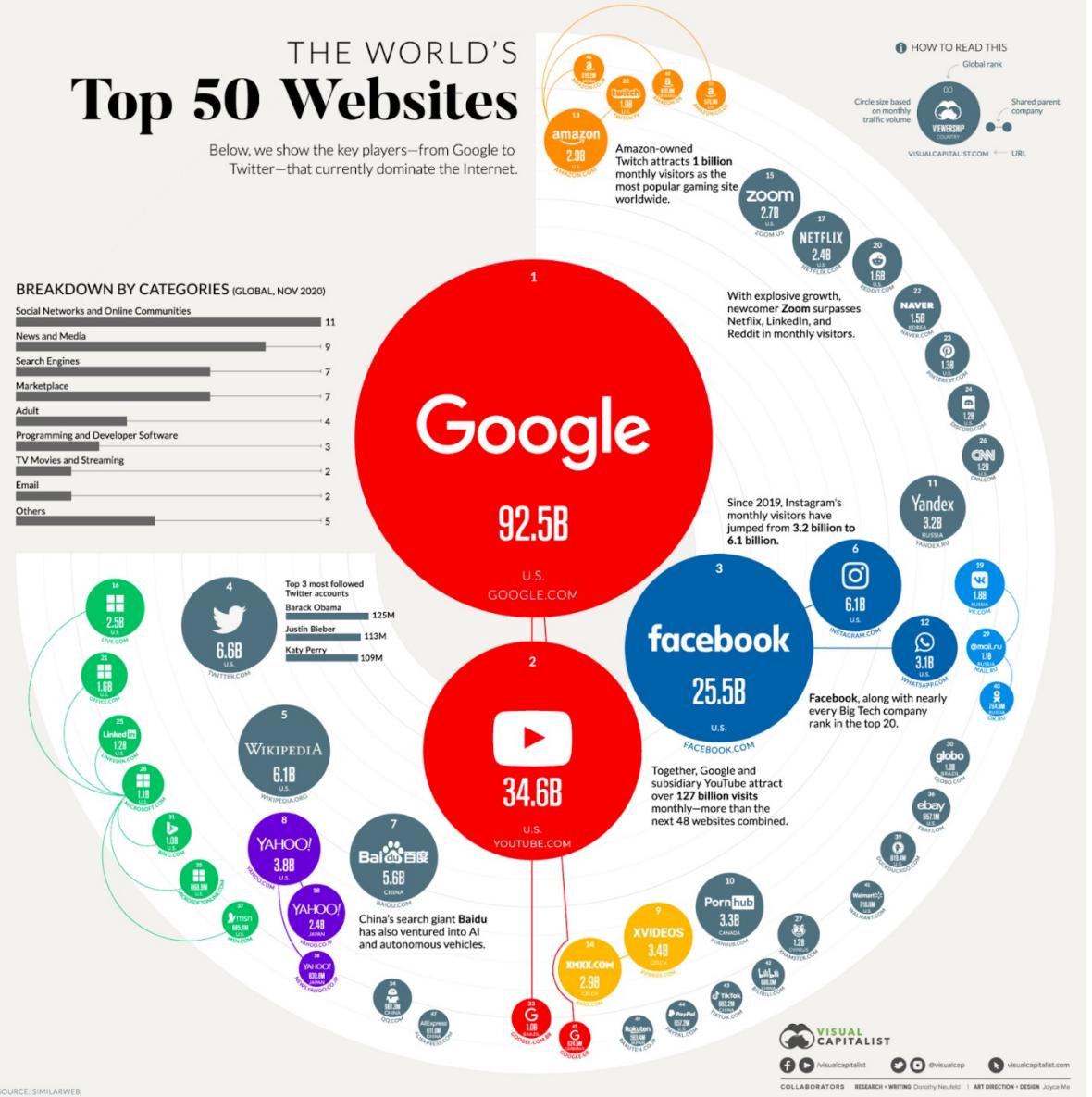
**Percentage of Billboard Hot 100 Number One Hits  
with a Key Change**



# Area comparison

# THE WORLD'S Top 50 Websites

Below, we show the key players—from Google to Twitter—that currently dominate the Internet.



# The Top 45 Richest celebrities in media/arts

Not including businessmen/women, inflation adjusted for deceased celebrities

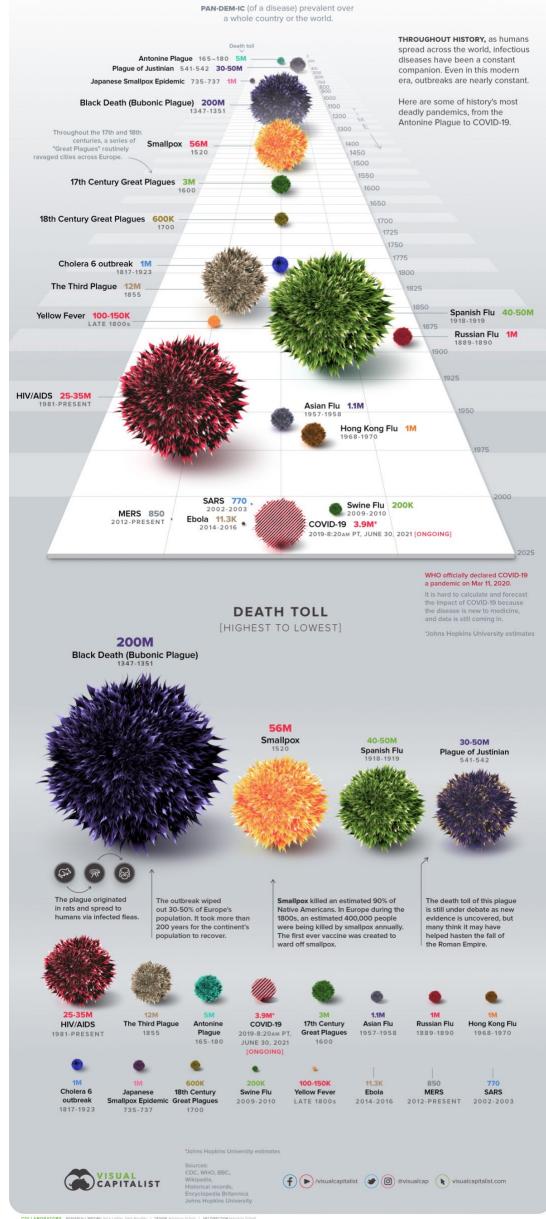


Source: [Celebritynetworth.com](http://Celebritynetworth.com)

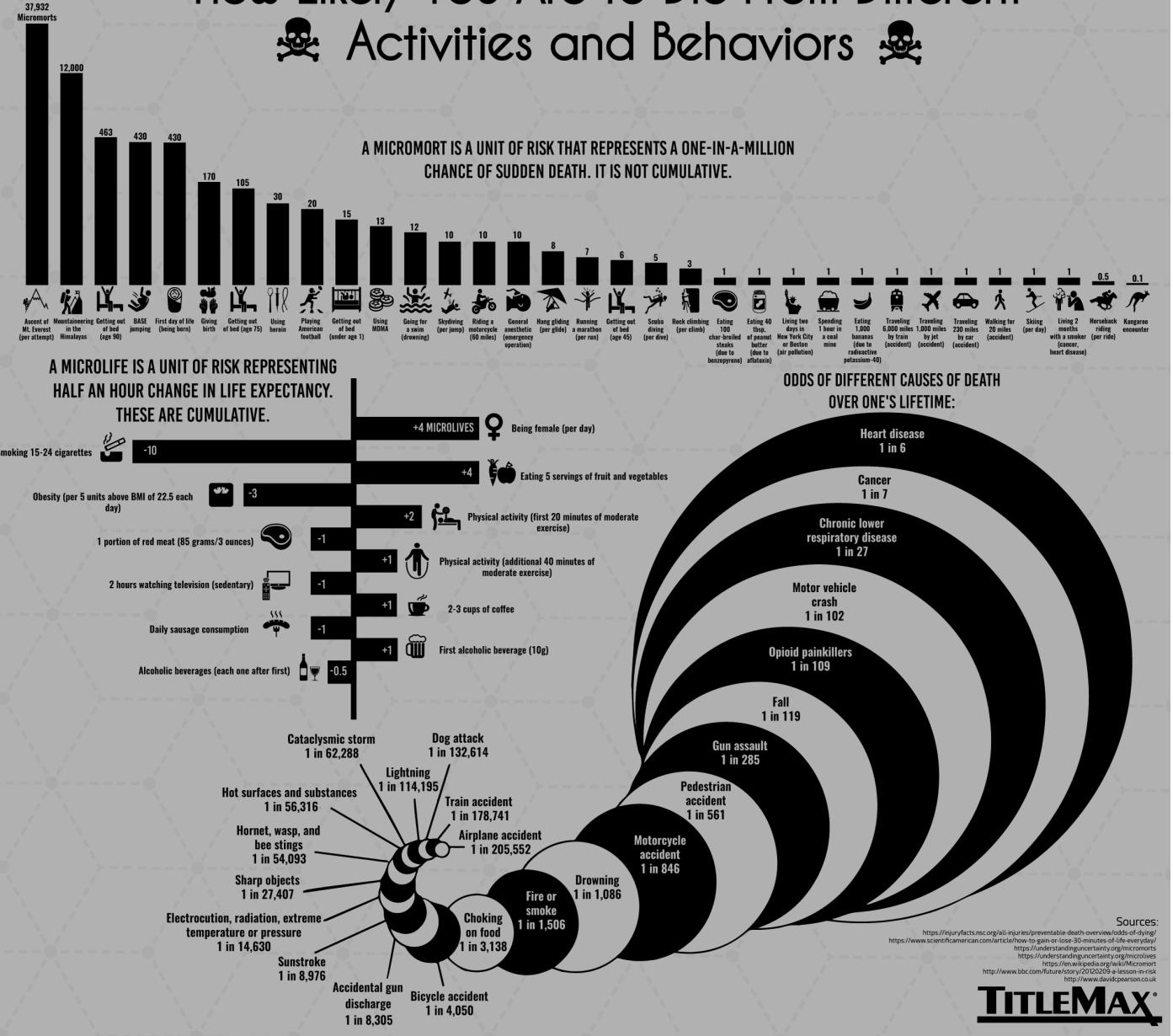
Created by  genuine  
impact

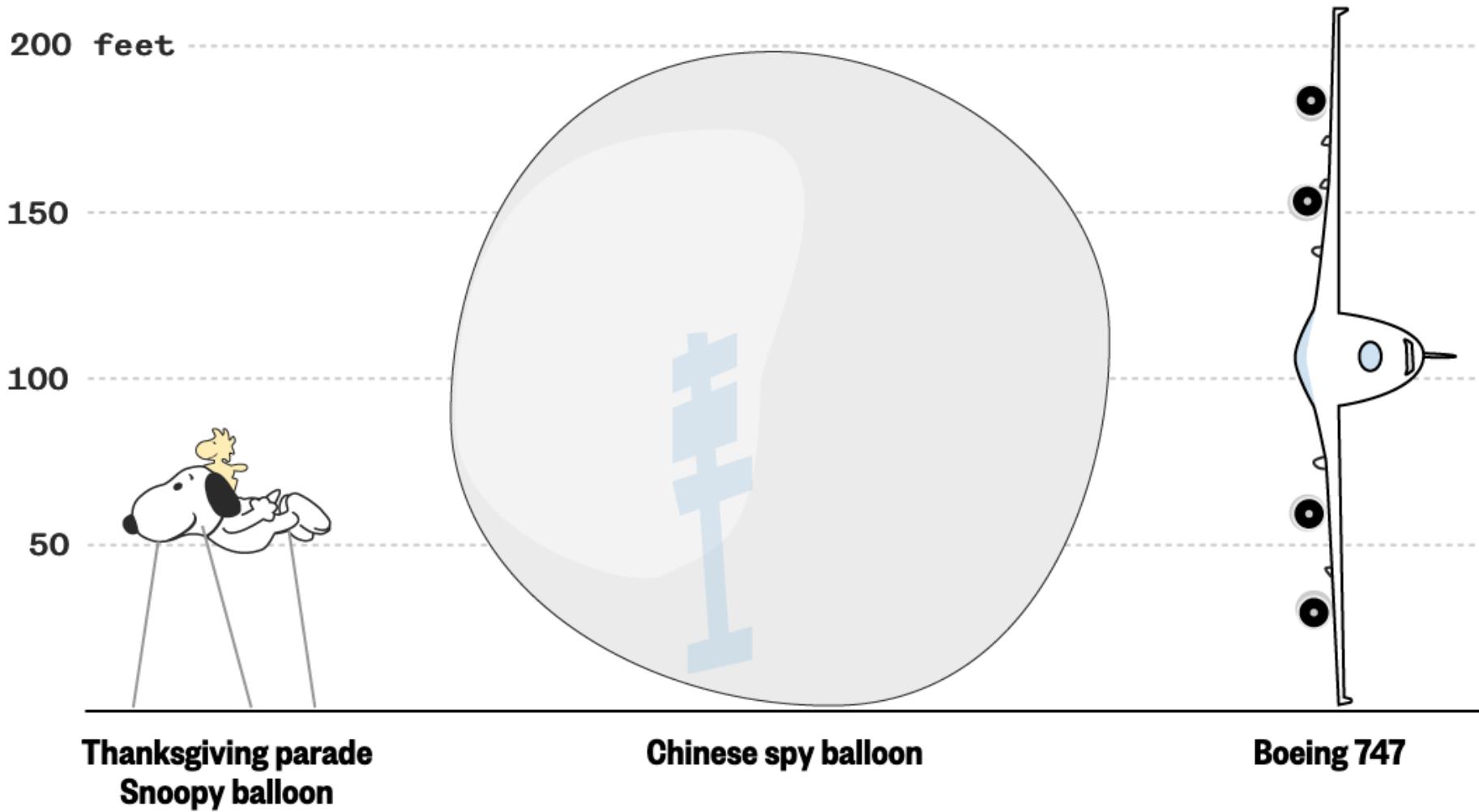
More charts: [genuineimpact.substack.com](https://genuineimpact.substack.com)

# HISTORY OF PANDEMICS



# How Likely You Are to Die From Different Activities and Behaviors

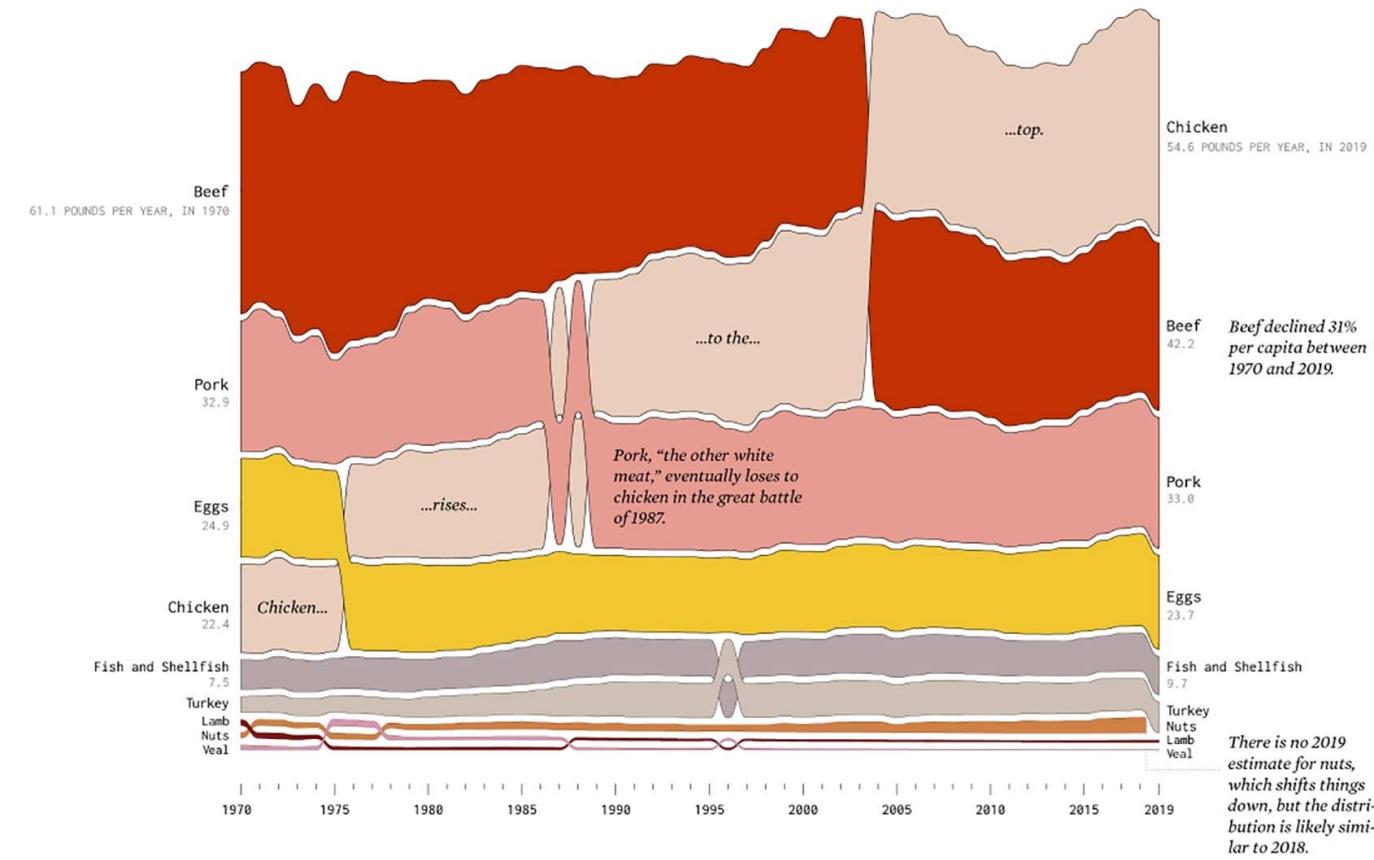




**Source:** NORAD, Boeing, NBC News

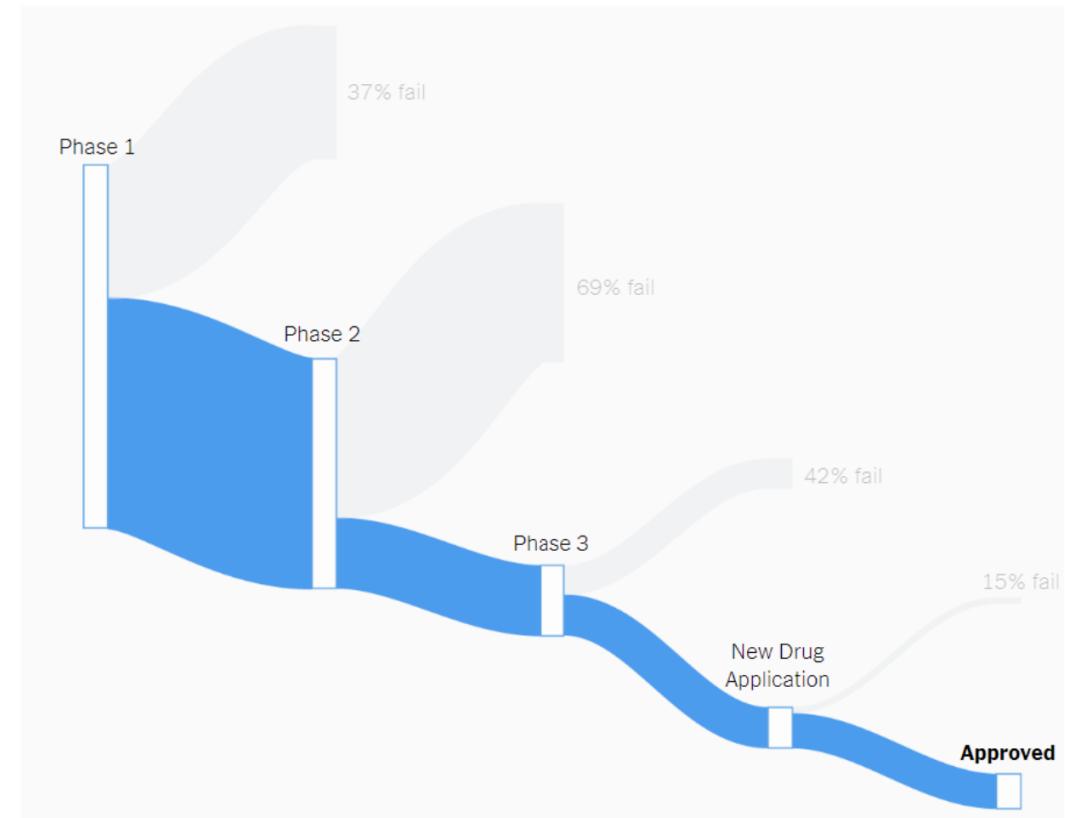
**Graphic:** JoElla Carman / NBC News

Stacked bargraphs/Sankey diagrams



## Less than 10 percent of drug trials are ultimately approved

Probability of success at each phase of research



# Ray tracing and animation



Data: GHS\_POP\_EU  
Software: Aerialod  
Author: Alasdair Rae

# Population Density in Europe

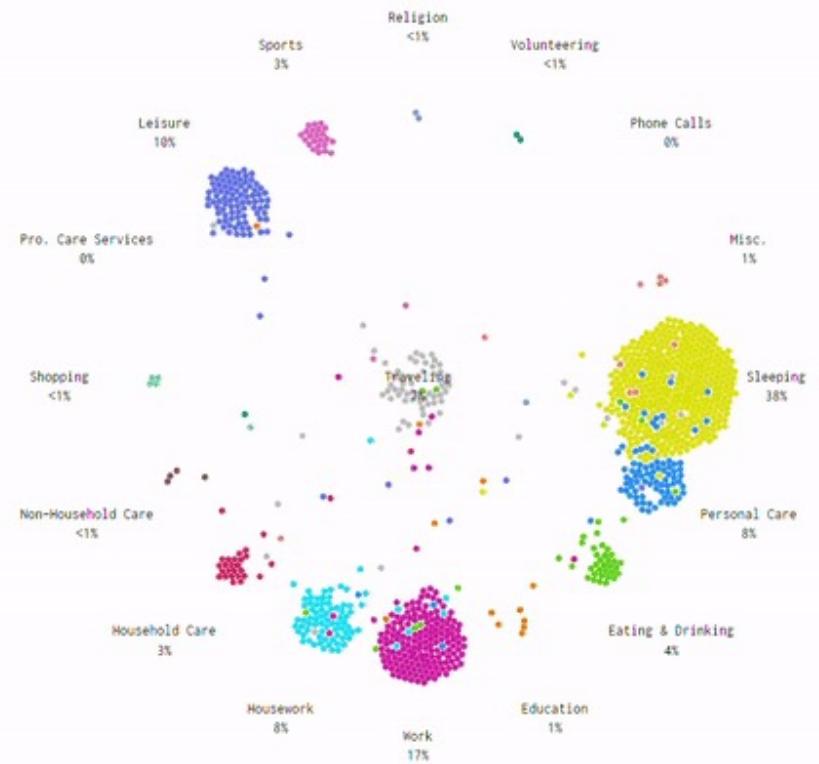
1km cells - the big squares are 50km x 50km

7:27am

SLOW MEDIUM FAST

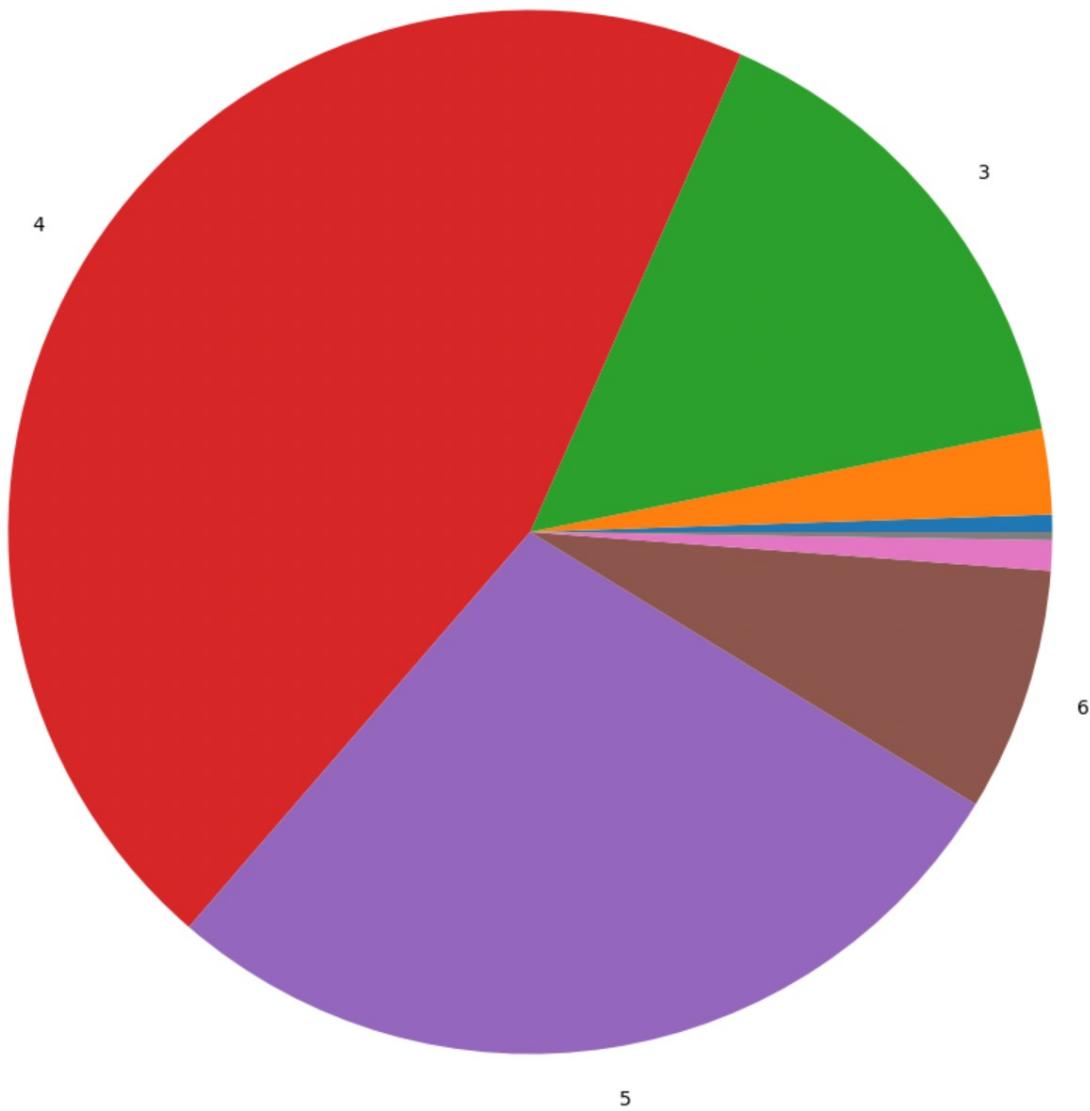
*It's wake up time for most. Time to start the day with morning rituals, breakfast and a wonderful commute.*

This is a simulation of 1,000 people's average day. It's based on 2014 data from the American Time Use Survey, made way more accessible by the ATUS Extract Builder.

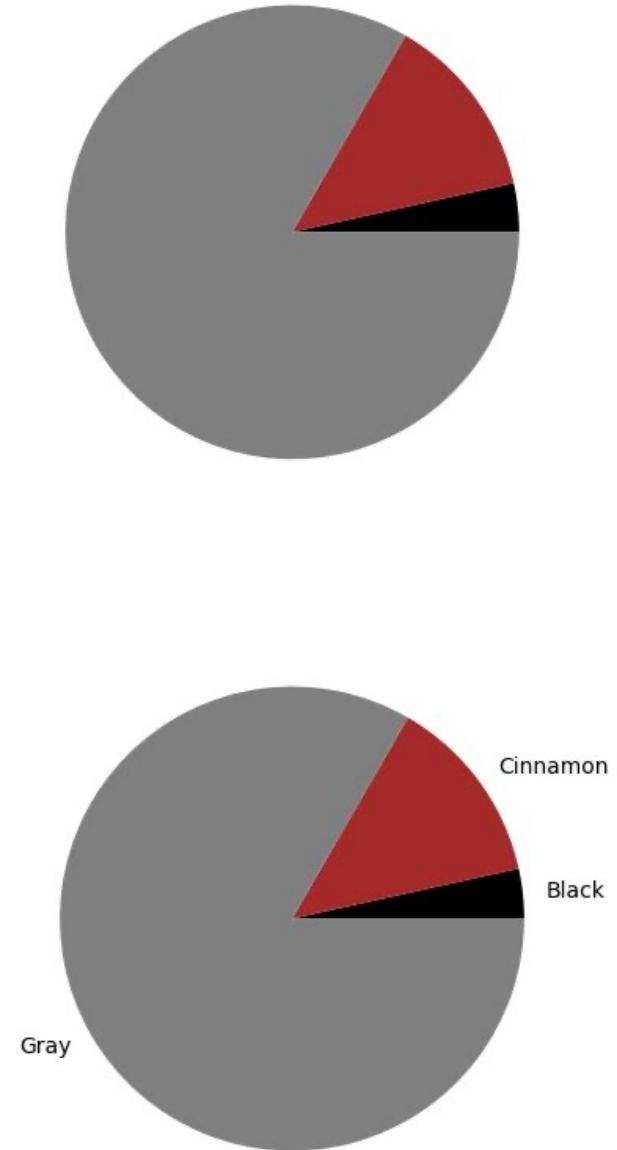


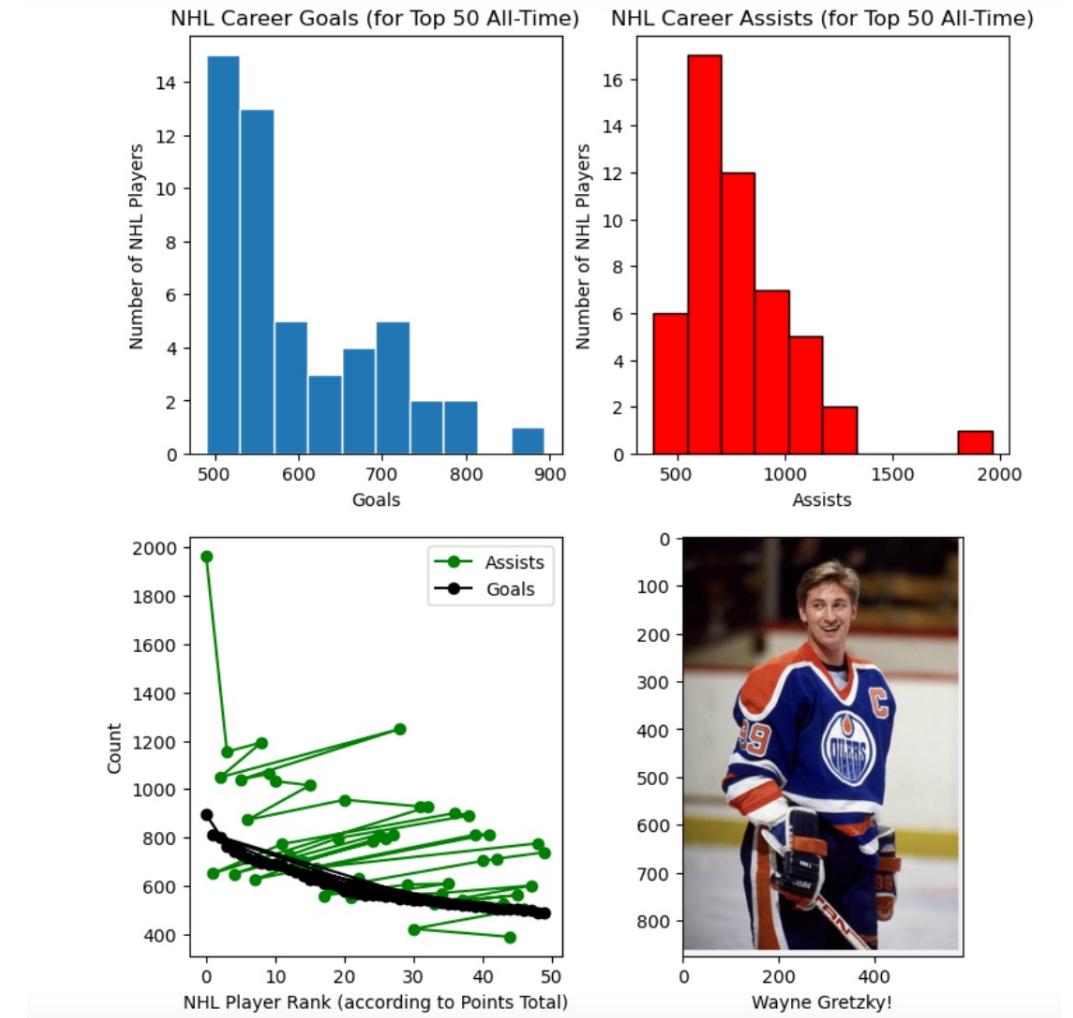
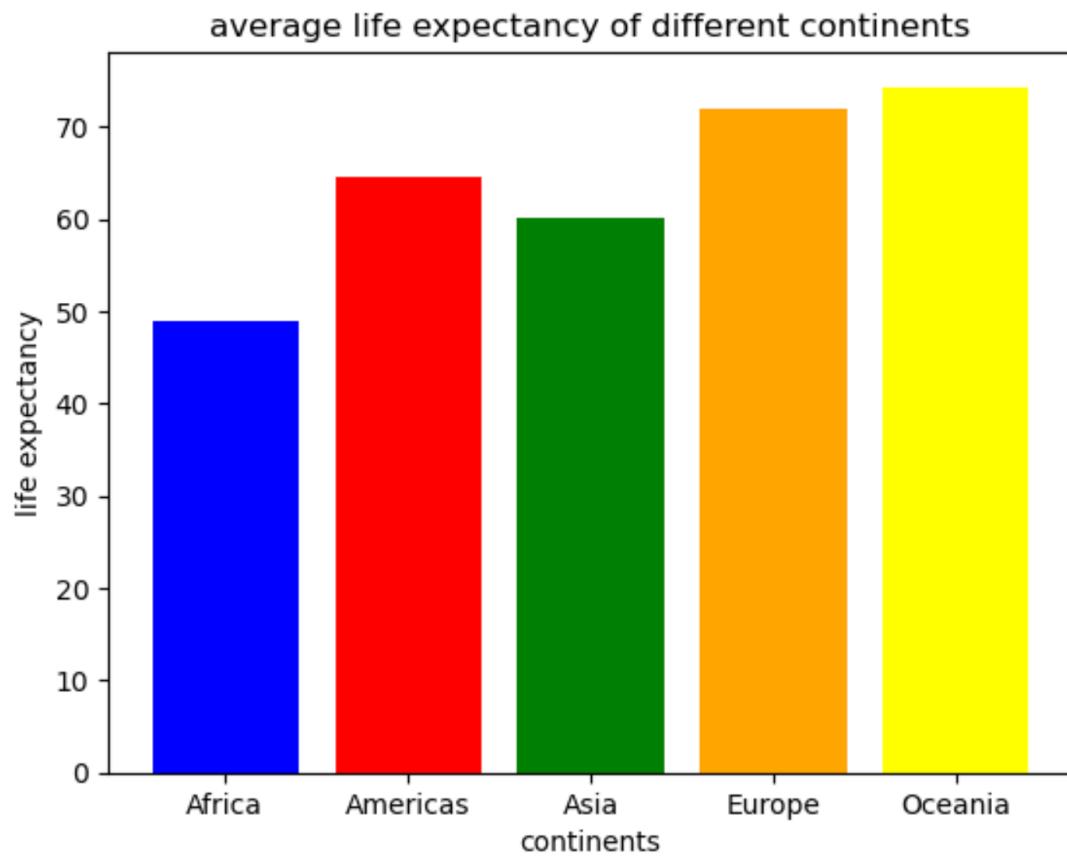
# Your visualizations

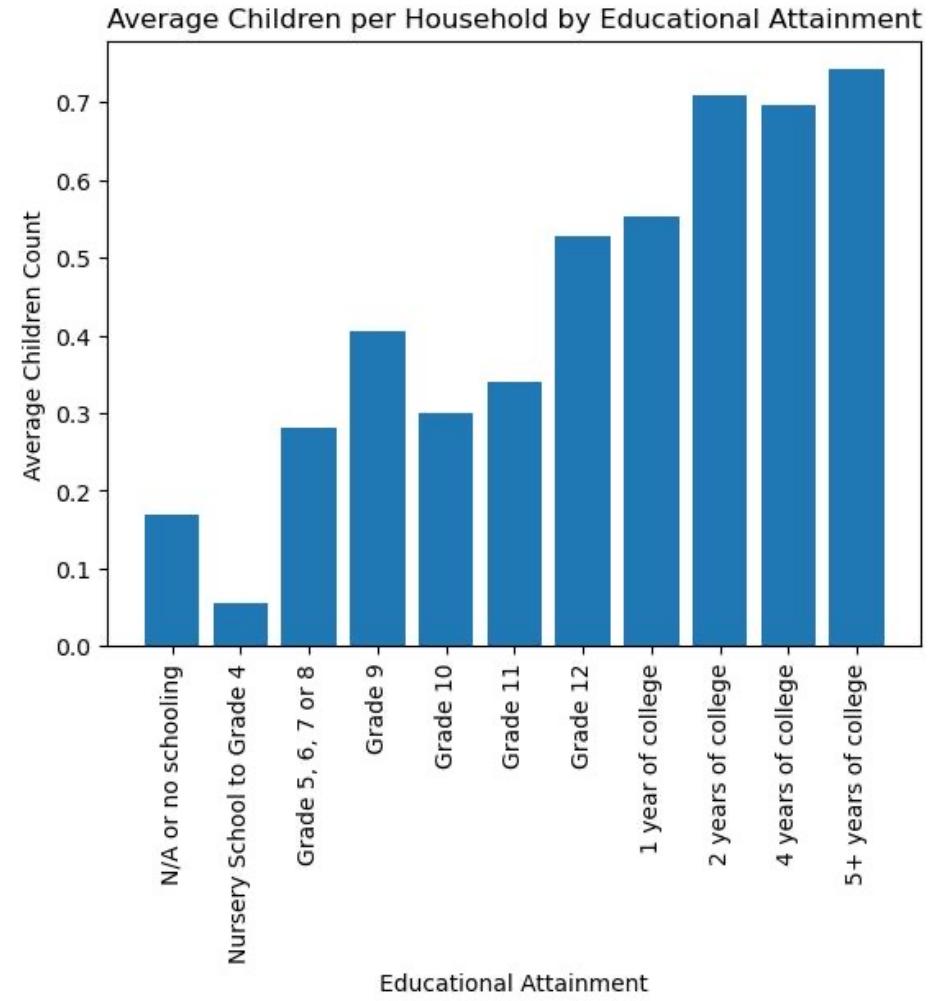
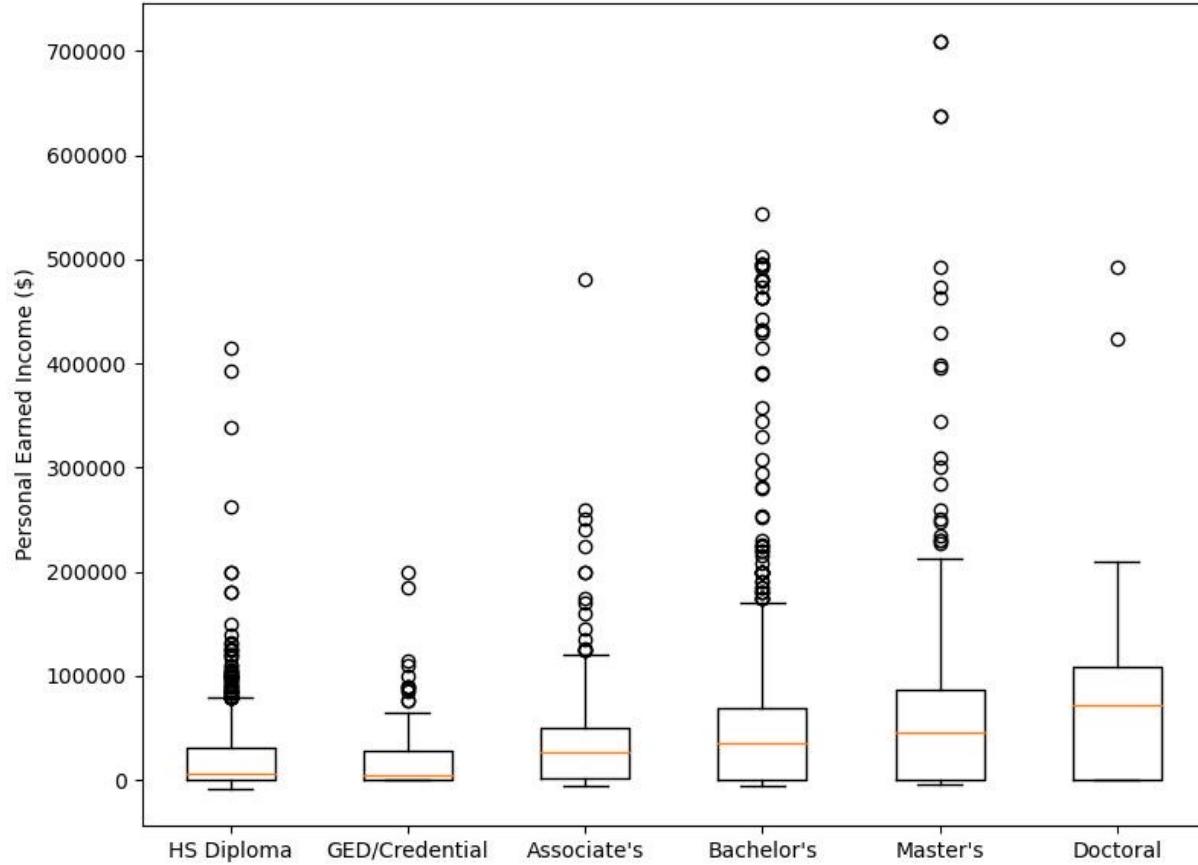
Number of Bedrooms in Each Household

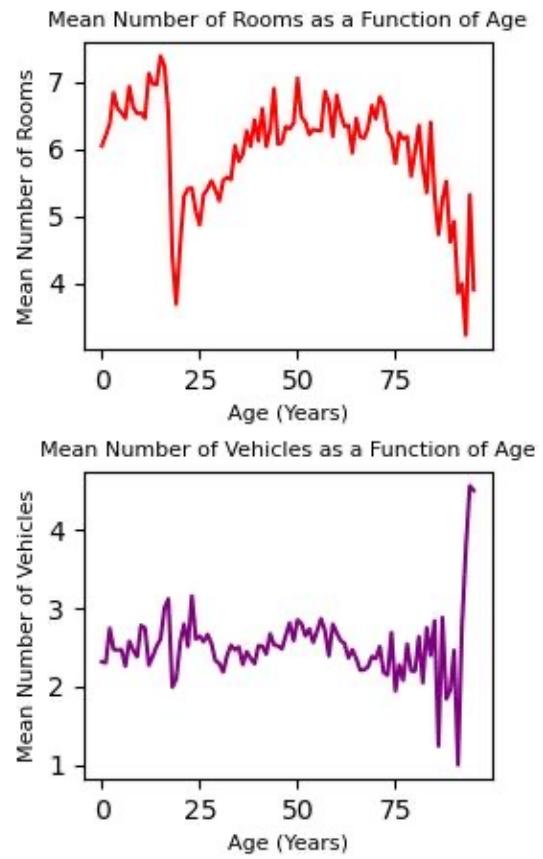
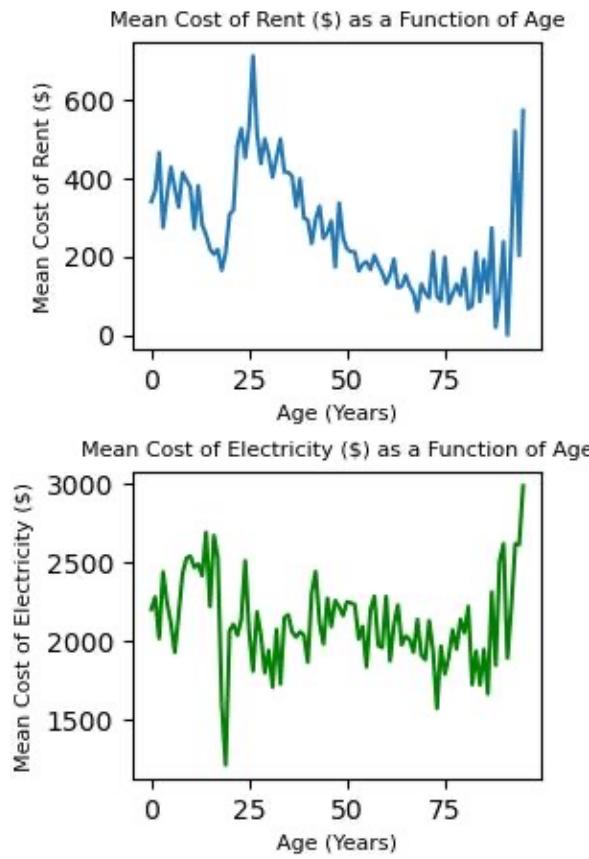
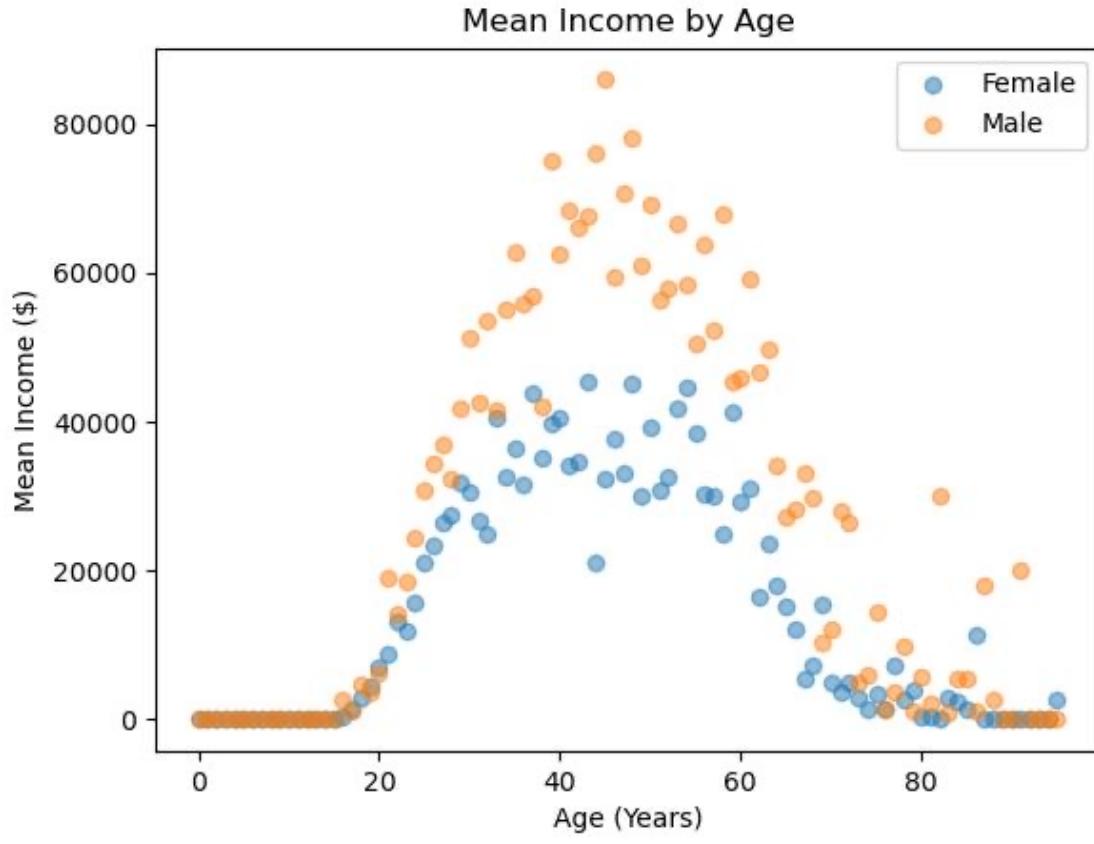


Indifferent Central Park Squirrels by Color

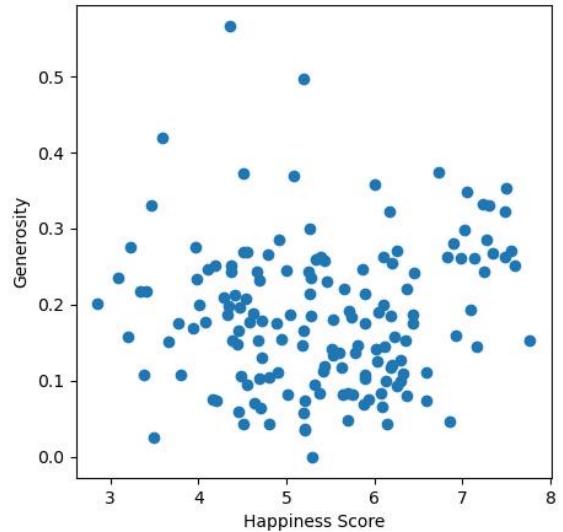




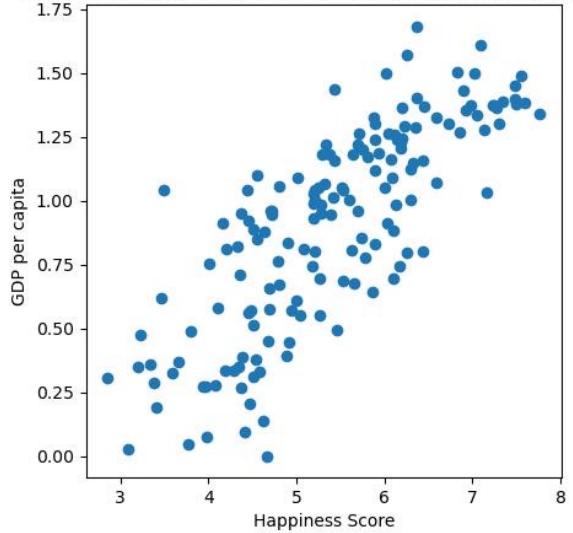




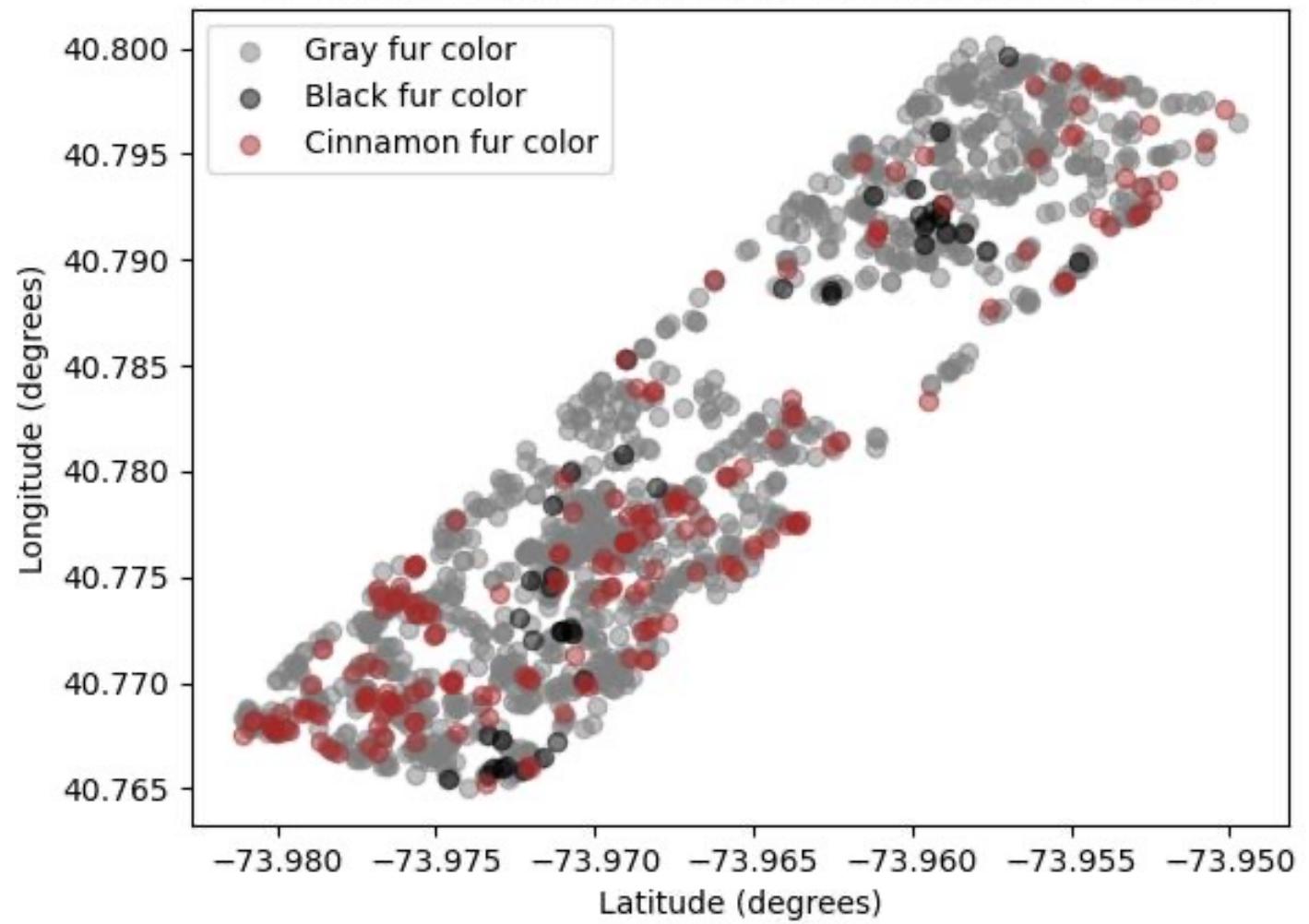
generosity as a function of a country's happiness score



GDP per capita as a function of a country's happiness score



Location of Indifferent Squirrels Based on Fur Color



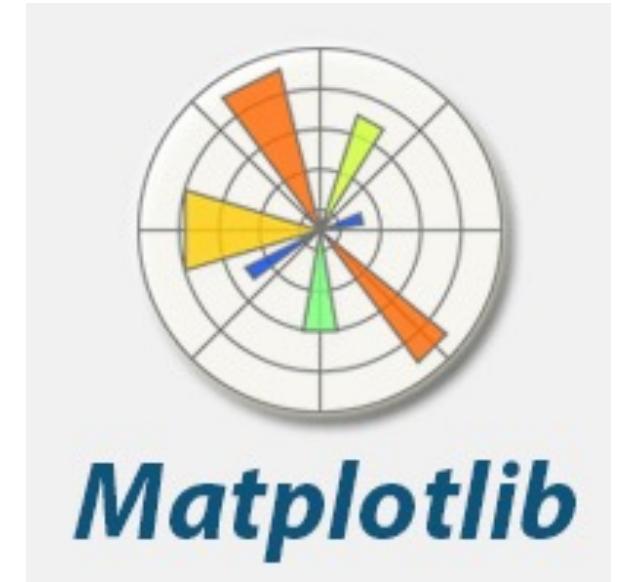
# Review: visualizing data with matplotlib

[Matplotlib](#) is a comprehensive library for creating static, animated, and interactive visualizations.

- `import matplotlib.pyplot as plt`

Types of plots we have created

- `plt.plot(x, y, '-o')` # line plot/scatter plot
- `plt.hist(data)`
- `plt.boxplot(data)`
- `plot.scatter(x, y, s = , color = , marker = )`



# Review: visualizing data with matplotlib

Make sure always label your axes:

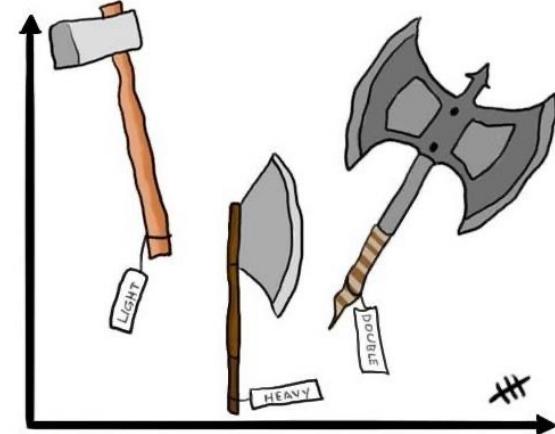
- `plt.ylabel("y label")`
- `plt.xlabel("x label")`
- `plt.title("my title")`
- `plt.plot(x, y, label = "blah")`
- `plt.legend()`

We can create subplots:

- `plt.subplot(1, 2, 1);`
- `plt.plot(x1, y1);`

Let's do some warm-up exercises Jupyter!

**Always label your axes**



# Seaborn

[“Seaborn](#) is a Python data visualization library based on [matplotlib](#). It provides a high-level interface for drawing attractive and informative statistical graphics.”

- i.e., it will create better looking plots that are easier to make



There are ways to create visualizations in seaborn:

1. **axes-level** functions that plot on a single axis
2. **figure-level** functions that plot across multiple axes

To make plots better looking we can set a theme

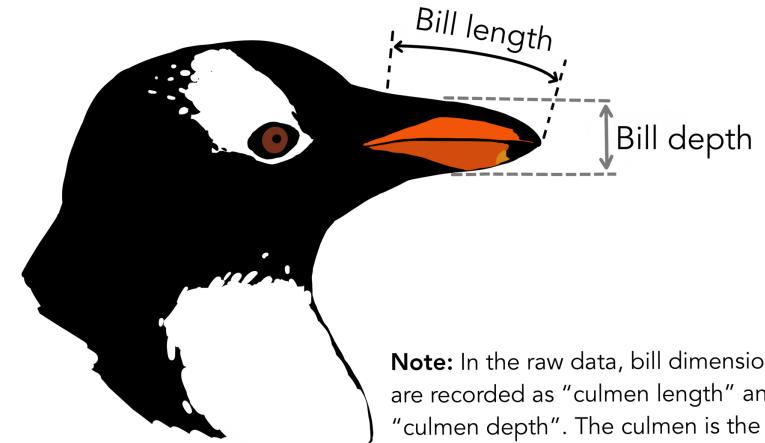
```
import seaborn as sns
```

```
sns.set_theme()
```

We will focus on figure level plots

# Inspiration: Palmer penguins

To explore seaborn, let's look at some data on penguins!



**Note:** In the raw data, bill dimensions are recorded as "culmen length" and "culmen depth". The culmen is the dorsal ridge atop the bill.

# Seaborn figure level plots

Figure level plots are grouped based on the types of variables being plotted

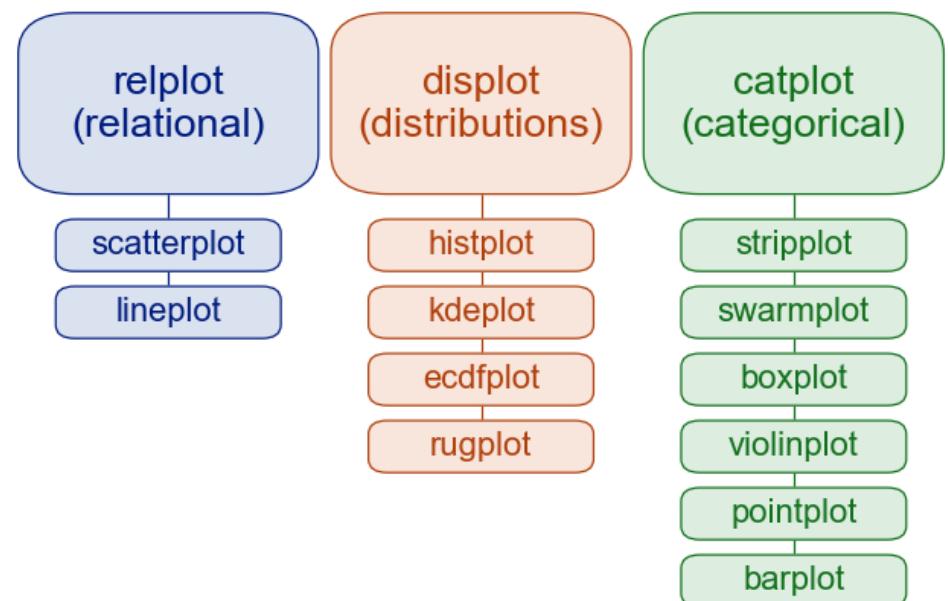
In particular, there are plots for:

1. Two quantitative variables
  - `sns.relplot()`
2. A single quantitative variable
  - `sns.displot()`

Quantitative variable compared across different categorical levels

- `sns.catplot()`

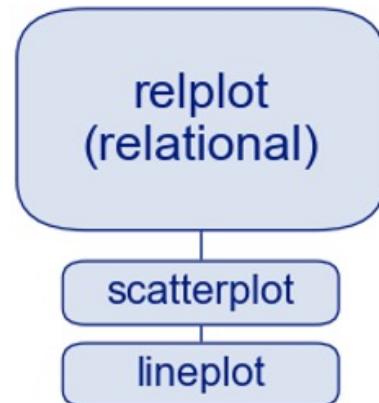
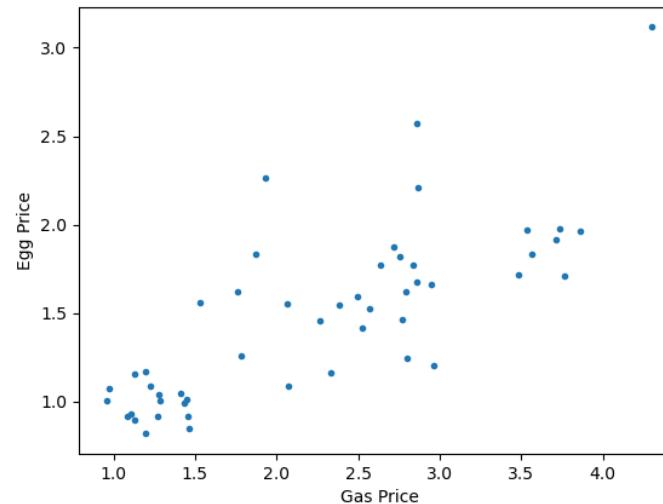
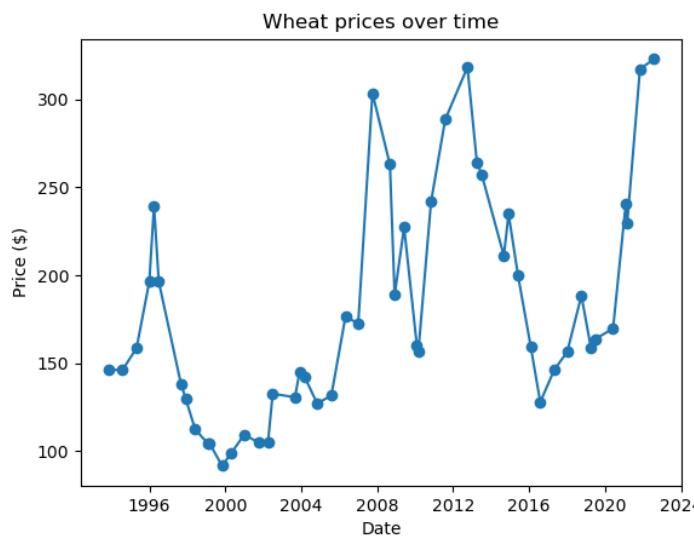
Figure level plots



# Plots for two quantitative variable

What types of plots have we seen for assessing the relationships between two quantitative variable?

- Line plots and scatter plots!

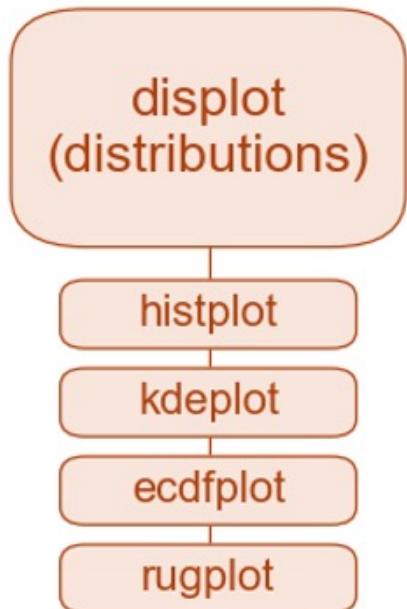
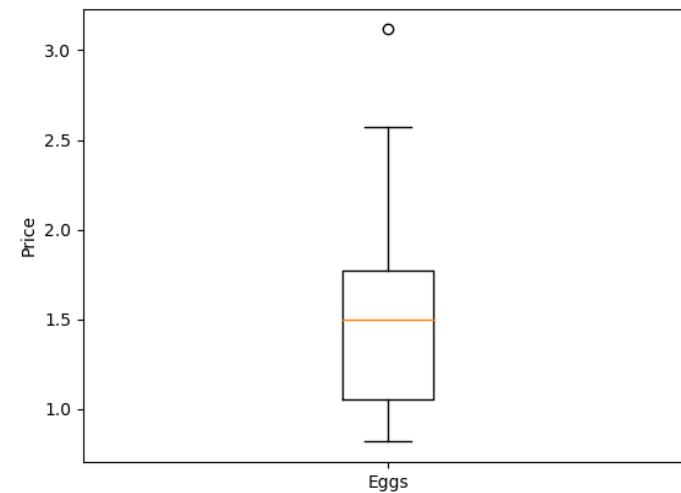
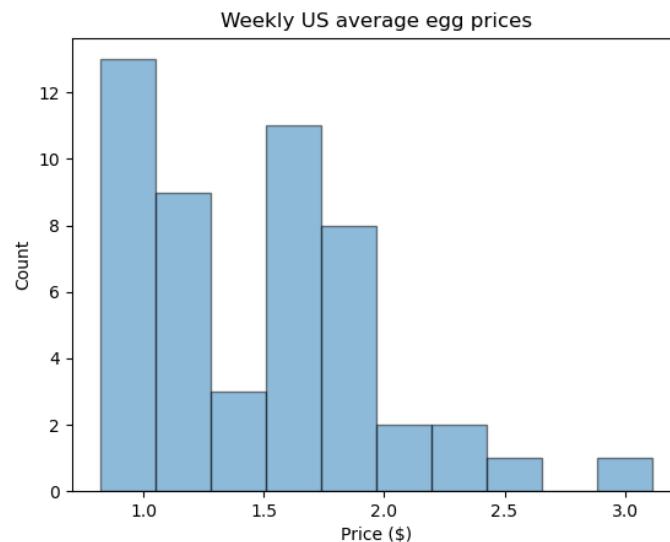


Let's explore this in Jupyter!

# Plots for a single quantitative variable

What types of plots have we seen for plotting a single quantitative variable?

- Histograms and boxplots

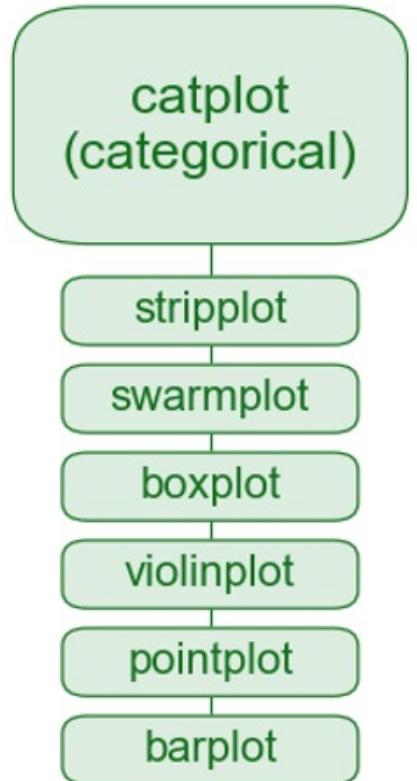
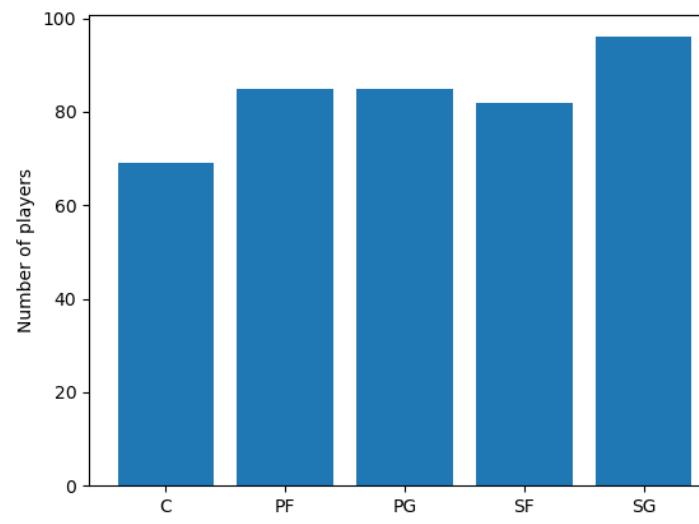
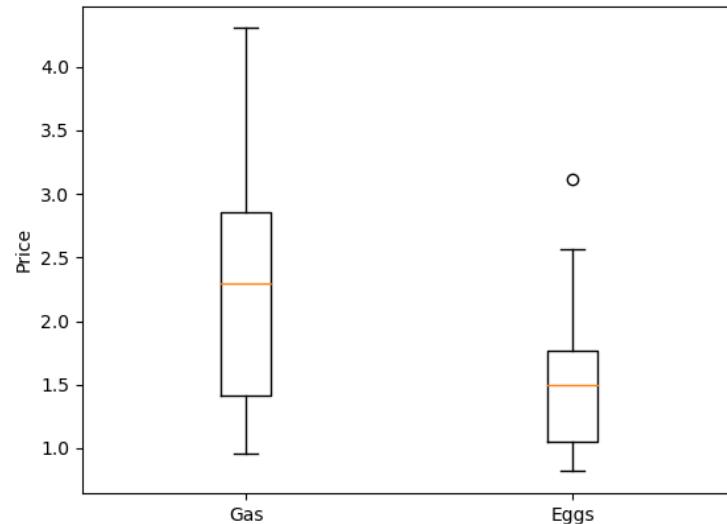


Let's explore this in Jupyter!

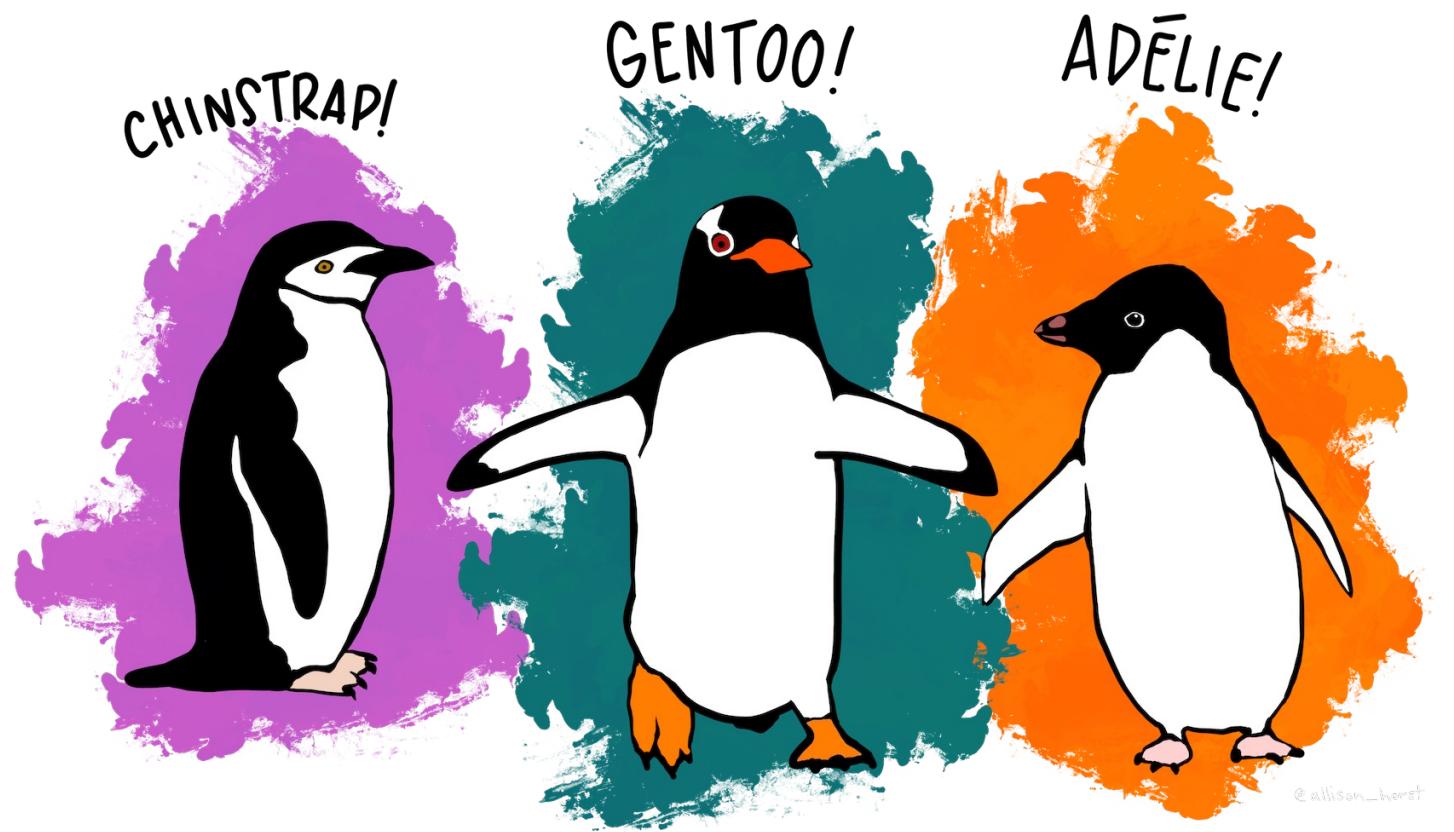
# Plots for quantitative data comparing across different categorical levels

What types of plots have we seen comparing quantitative data at different levels of a categorical variable?

- Side-by-side boxplots, barplots (sort of)



Let's explore this in Jupyter!



@allison\_horst

**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

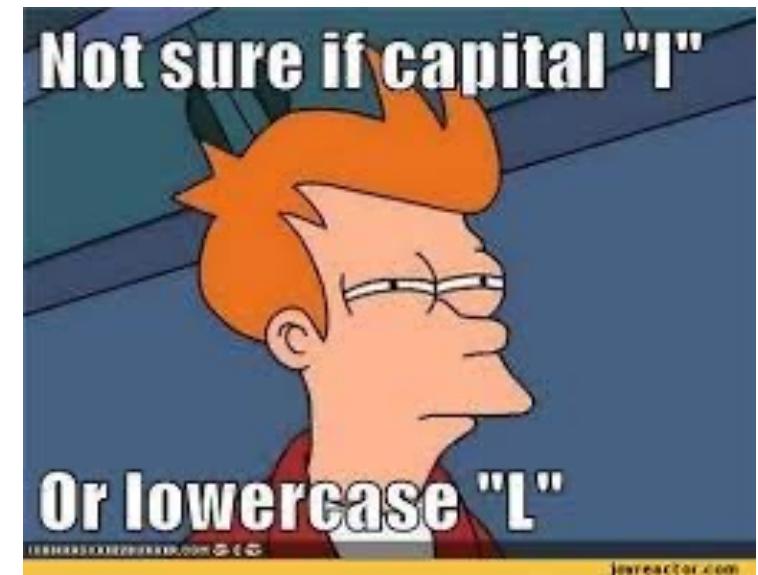
Python has many string methods that are useful for manipulating text and cleaning data!

# Text manipulation: capitalization

Some of the simplest string methods involve changing capitalization.

Changing capitalization can be useful when joining DataFrames

- i.e., if their 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!