# YData: Introduction to Data Science



# Class 11: Intro to data visualization

# Overview

Quick review of pandas DataFrames and joining DataFrames

History of Data Visualization

Visualizing data with matplotlib!

# Announcement: Homework 5

Homework 5 has been posted!

It is due on Gradescope on <span style="color:red">Sunday February 26th at 11pm</span>

- **Be sure to mark each question on Gradescope along with the <span style="color:red">page that has the answers</span>!**

# Announcement: 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

- It's a chance to practice everything you've learned in class!

The goal of is to present a clear and compelling data analysis showing a few interesting results!

- Ideally something you are proud of that you could show off to potential future employers, etc.

A few sources for data sets are listed on Canvas

- You can use data you collect as well. If you use data for another class your work must be unique for each class.

# Announcement: class project

A draft of the project is due on April 7<sup>th</sup>
- I'm telling you about this early so:
  - If you want to think/work on the project over break you can
  - If don't want to think about it over break you don't have to

There will be a "peer review" period where you will give and receive feedback on three other projects
- Instructions for how to do the review will be given

The final version will be due on April 30<sup>th</sup>

A Jupyter notebook "template" project will be available soon


ALL YOU NEED IS MOTIVATION

FALSE: YOU NEED FEAR AND AN APPROACHING DEADLINE
memegenerator.net

Questions?

# Quick review: pandas DataFrames

| PLAYER | POSITION | TEAM | SALARY |
|---|---|---|---|
| str | str | str | f64 |
| "Paul Millsap" | "PF" | "Atlanta Hawks" | 18.671659 |
| "Al Horford" | "C" | "Atlanta Hawks" | 12.0 |
| "Tiago Splitter... | "C" | "Atlanta Hawks" | 9.75625 |
| "Jeff Teague" | "PG" | "Atlanta Hawks" | 8.0 |
| "Kyle Korver" | "SG" | "Atlanta Hawks" | 5.746479 |

Pandas DataFrame hold Table data

Selecting columns:
- my_df[["col1", "col2"]].copy()      # getting multiple columns using a list

Extracting rows:
- my_df.iloc[0]                       # getting a row by number
- my_df.loc["index_name"]             # getting a row by Index value
- my_df [my_df["col_name"] == 7]   # getting rows using a Boolean mask

# Quick review: pandas DataFrames

Sorting rows of a DataFrame

    `my_df.sort_values(“col_name”, ascending = False)`   `# sort from largest to smallest`

Adding a new:

- `my_df[“new_col”] = values_array`

Renaming a column:

- `rename_dictionary = {”old_col_name": ”new_col_name"}`
- `my_df.rename(columns = rename_dictionary )`

# Creating aggregate statistics by group

We can get statistics separately by group:

- dow.groupby("Year").agg("max")



my_df.groupby("group_col_name").agg(

    new_col1 = ('col_name', 'statistic_name1'),

    new_col2 = ('col_name', 'statistic_name2'),

    new_col3 = ('col_name', 'statistic_name3')

)

# Creating a DataFrame "by hand"

We can create small DataFrames by entering data manually from a dictionary using the df.DataFrame() function

To do this we first create a dictionary, where:

- The **keys** are the names of the columns
- The **values** are the data in the rows

my_dictionary = {"col1": ["a", "b", "c"], "col2": [1, 2, 3]}

my_df = pd.DataFrame(my_dictionary)

# Review: Joining data frames

Suppose we have two DataFrames (or Series) called **x_df** and **y_df**

- x_df have two columns called key_x, and  val_x
- y_df has two columns called key_y and val_y

key_x          val_x

| 1 | x1 |
| 2 | x2 |
| 3 | x3 |

**DataFame: x_df**
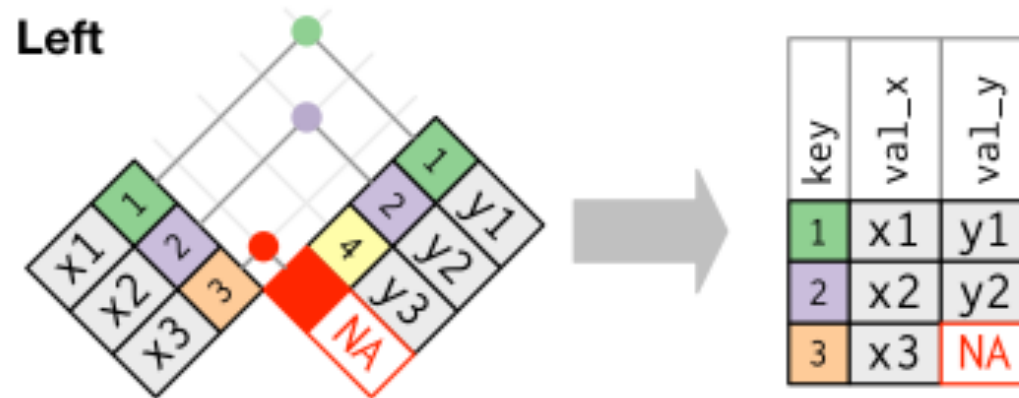
key_y          val_y

| 1 | y1 |
| 2 | y2 |
| 4 | y3 |

**DataFrame y_df**

Joins have the general form:

x_df.merge(y_df, left_on = "key_x",  right_on = "key_y")

# Left joins

**Left joins** keep all rows in the <u>left</u> table.

Data from <u>right</u> table is added when there is a matching key, otherwise NA as added.



x_df.merge(y_df, how = "left", left_on = "key_x",  right_on = "key_y")

# Right joins

**Right joins** keep all rows in the <u>right</u> table.

Data from <u>left</u> table added when there is a matching key, otherwise NA as added.



x_df.merge(y_df, how = "right", left_on = "key_x",  right_on = "key_y")

# Inner joins

**Inner joins** only keep rows in which there are matches between the keys in both tables.



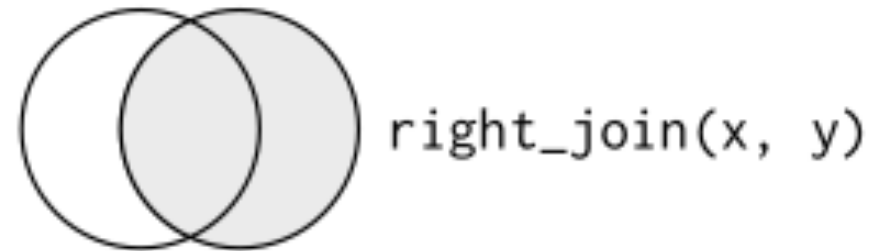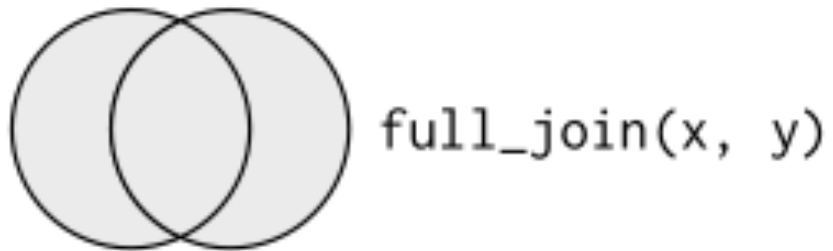x_df.merge(y_df, how = "inner", left_on = "key_x",  right_on = "key_y")
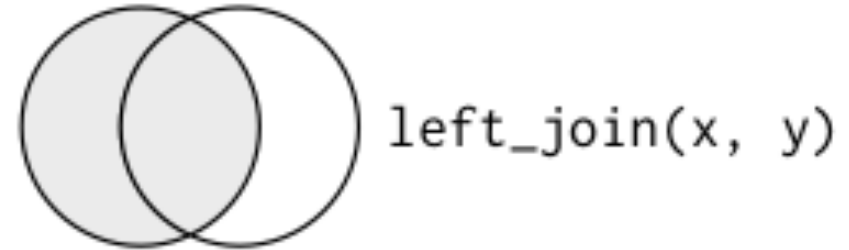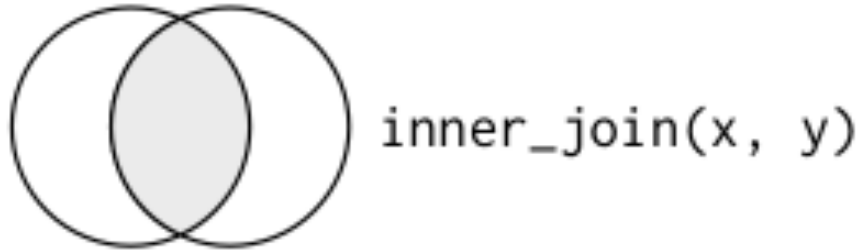
# Full (outer) joins

**Full joins** keep all rows in both table.

NAs are added where there are no matches.



x_df.merge(y_df, how = "outer", left_on = "key_x",  right_on = "key_y")

# Summary

# Joining on Index values

If two DataFrames **have the same Index values** then we can join them using the .join() method instead of the .merge() method

The .join() method is very similar to .merge() except we don't need to specify left_on and right_on arguments since the DataFrames are being joined by their Indexes

An example of a left join would be:
- x_df.join(y_df, how = "left")    # assuming x_df and y_df has the same Index values

Let's review this in Jupyter!

# Questions?

# Data visualization!

# A very brief history of data visualization...

## The Golden Age of Statistical Graphics

**Michael Friendly**

# Data visualization

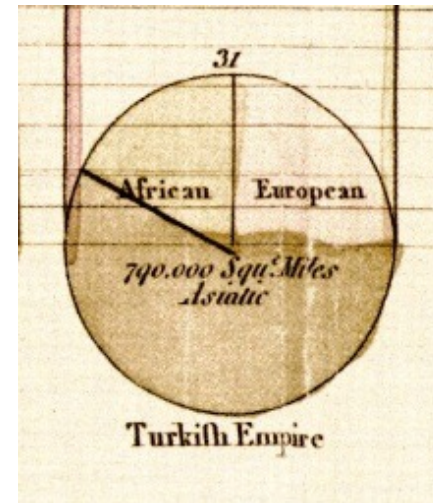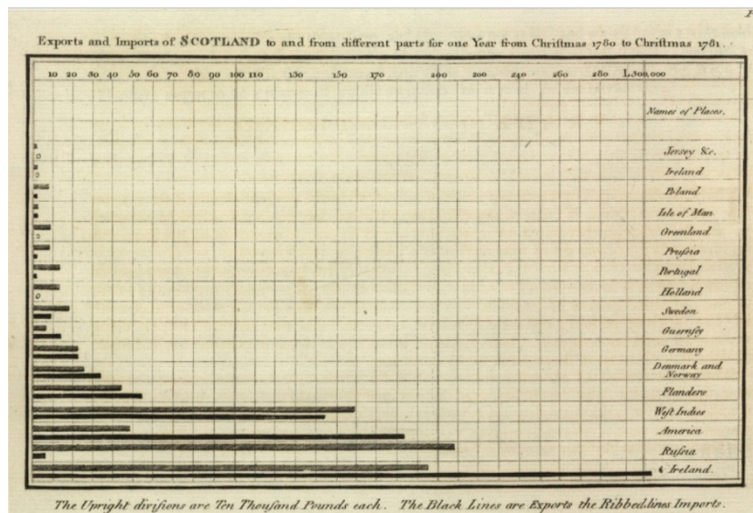What are some reasons we visualize data rather than just reporting statistics?

*Whatever relates to extent and quantity may be represented by geometrical figures. 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*
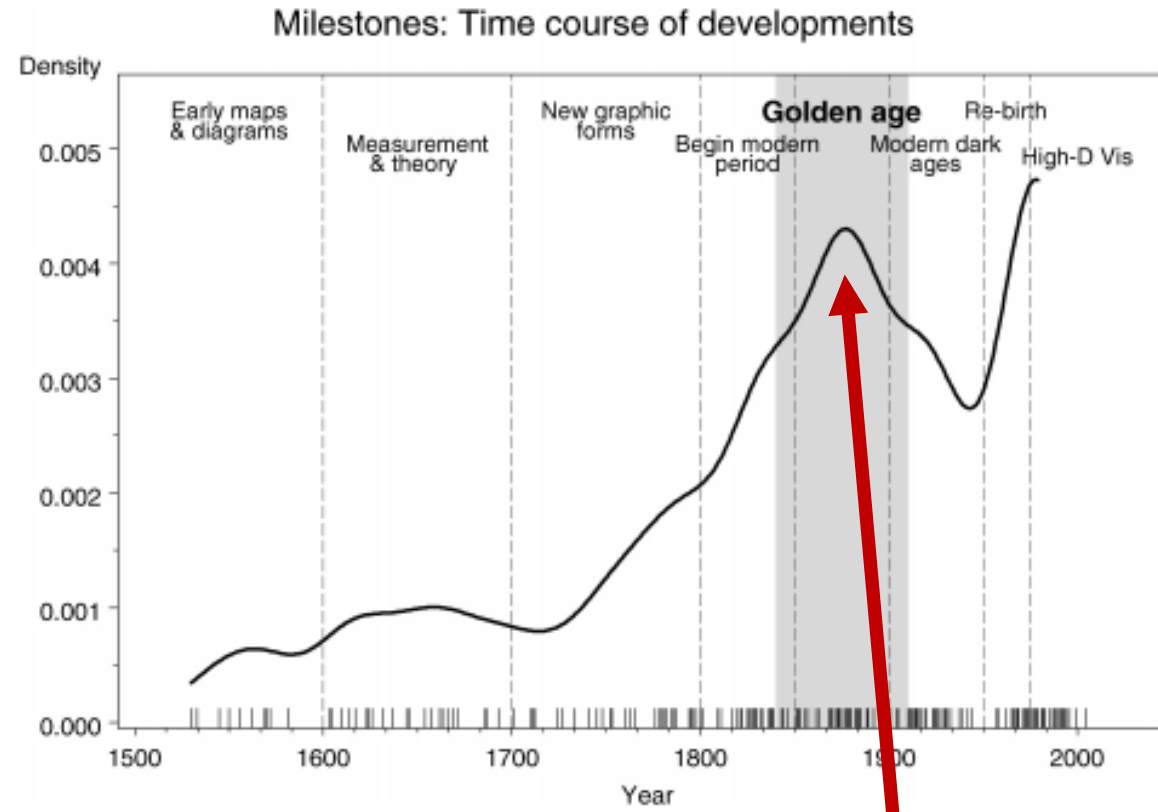
# A very brief history of data visualization

The age of modern statistical graphs began around the beginning of the 19th century

William Playfair (1759-1823) is credited with inventing the line graph, bar chart and pie chart
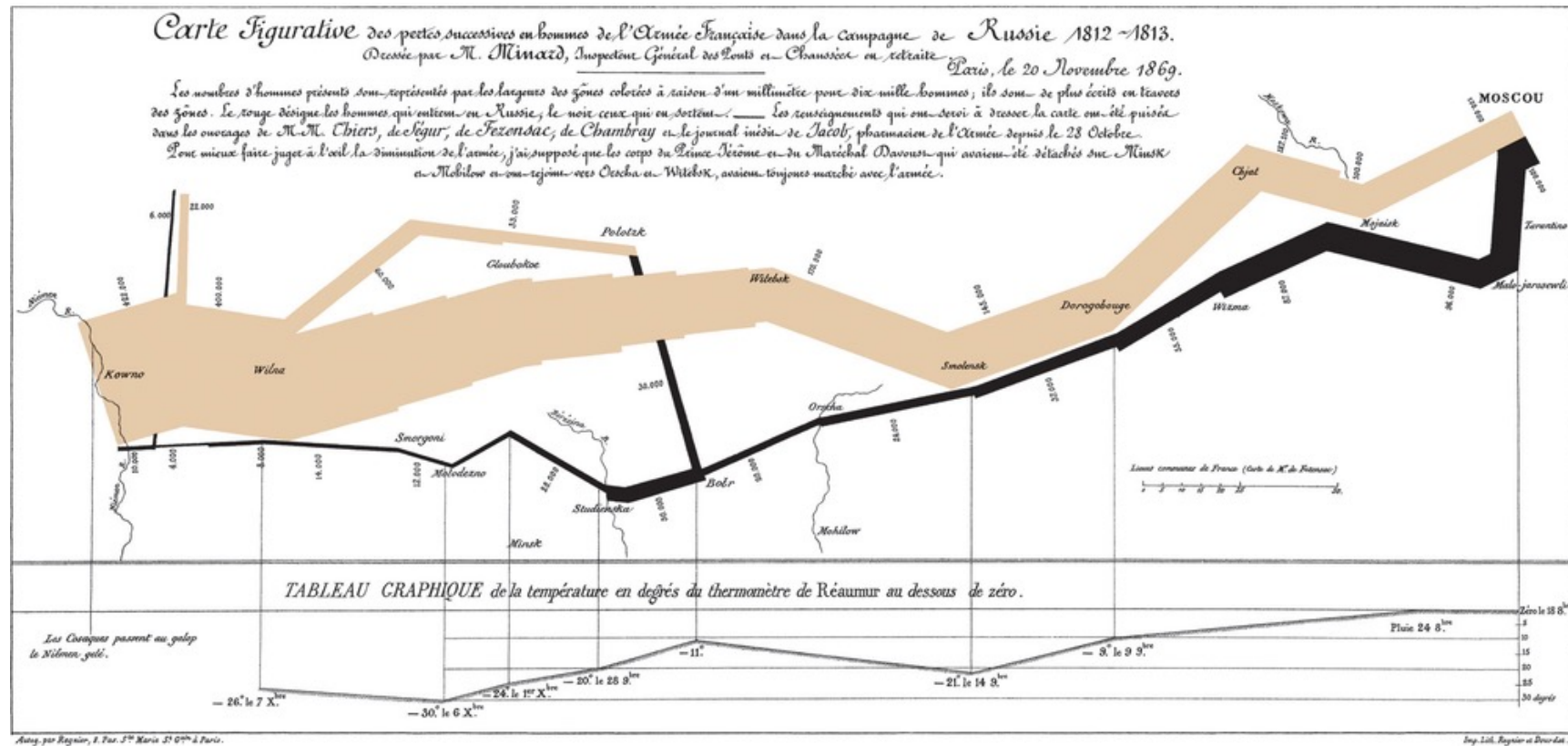
# A very brief history of data visualization

According to Friendly, statistical graphics researched its golden age between 1850-1900

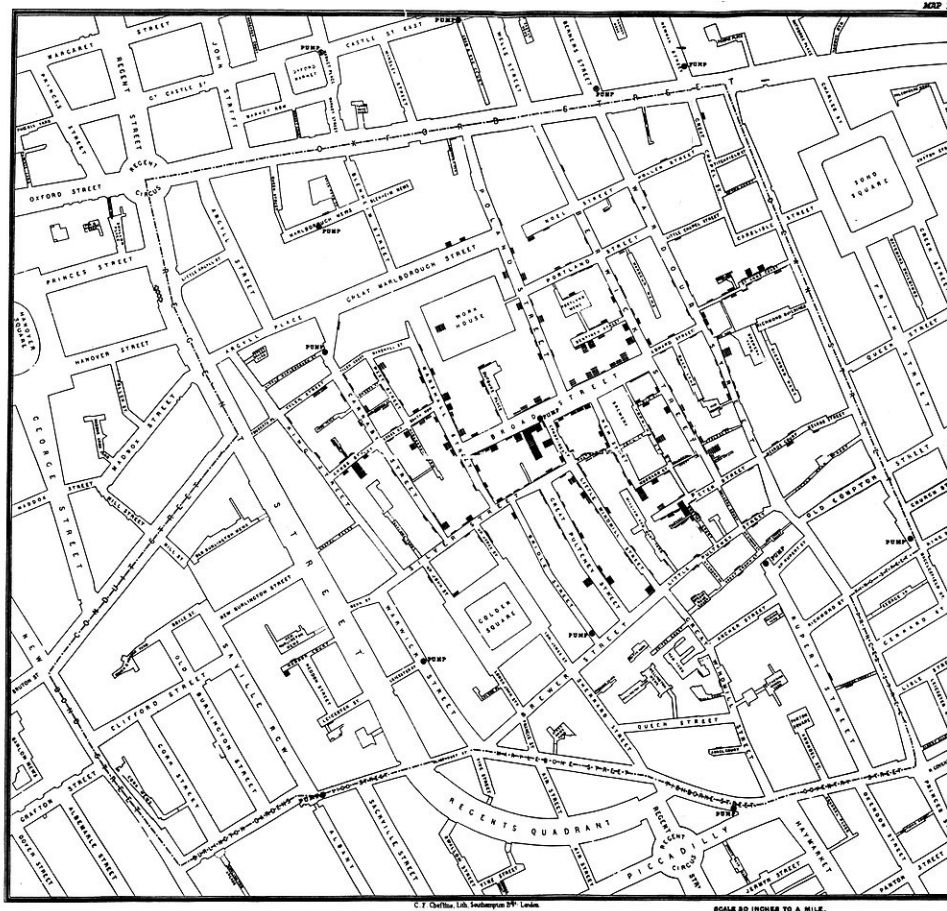# A very brief history of data visualization

[Joseph Minard](#) (1781-1870)



Map of Napoleon's march on Russia

# A very brief history of data visualization

John Snow (1813-1858)



Clusters of cholera cases in London epidemic of 1854

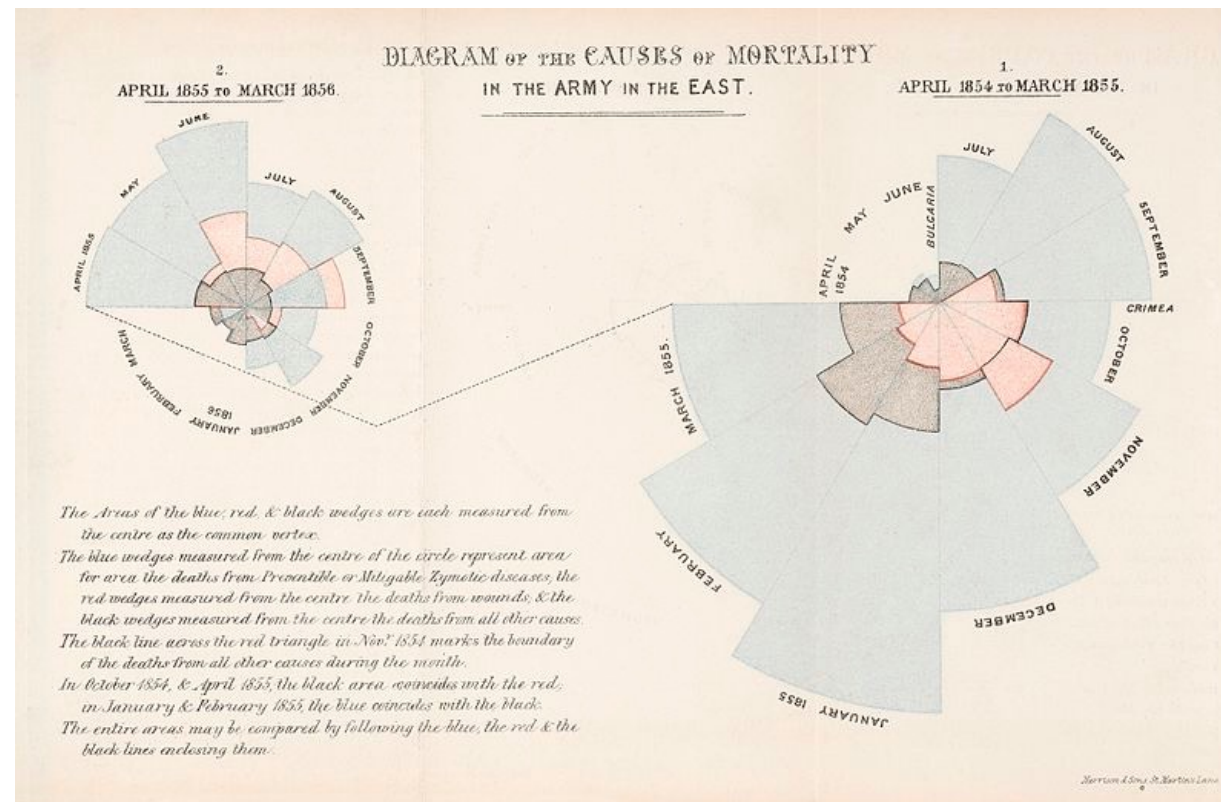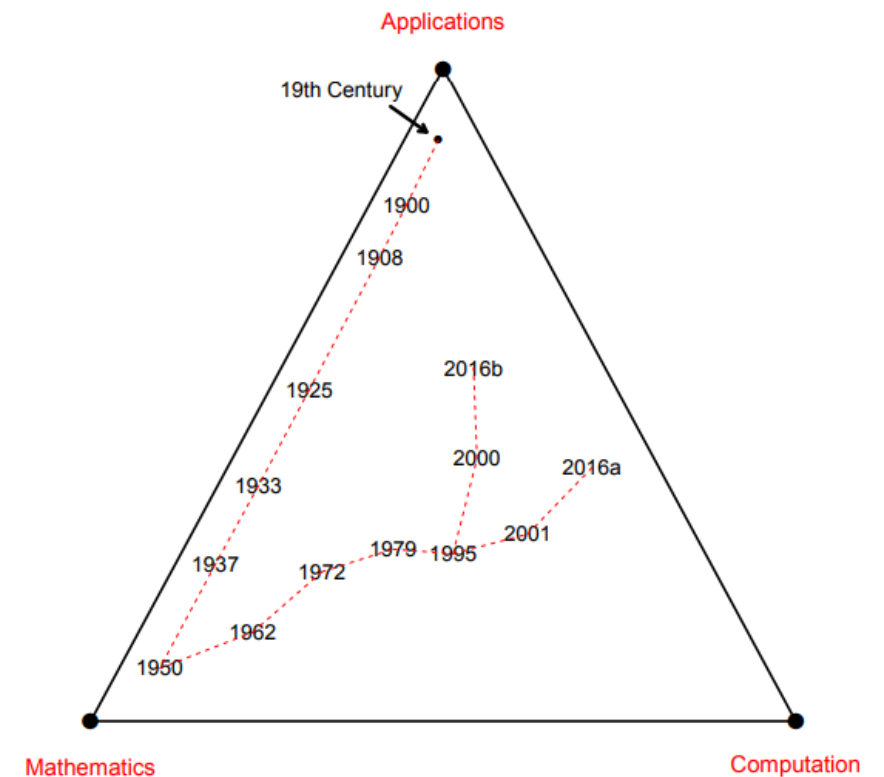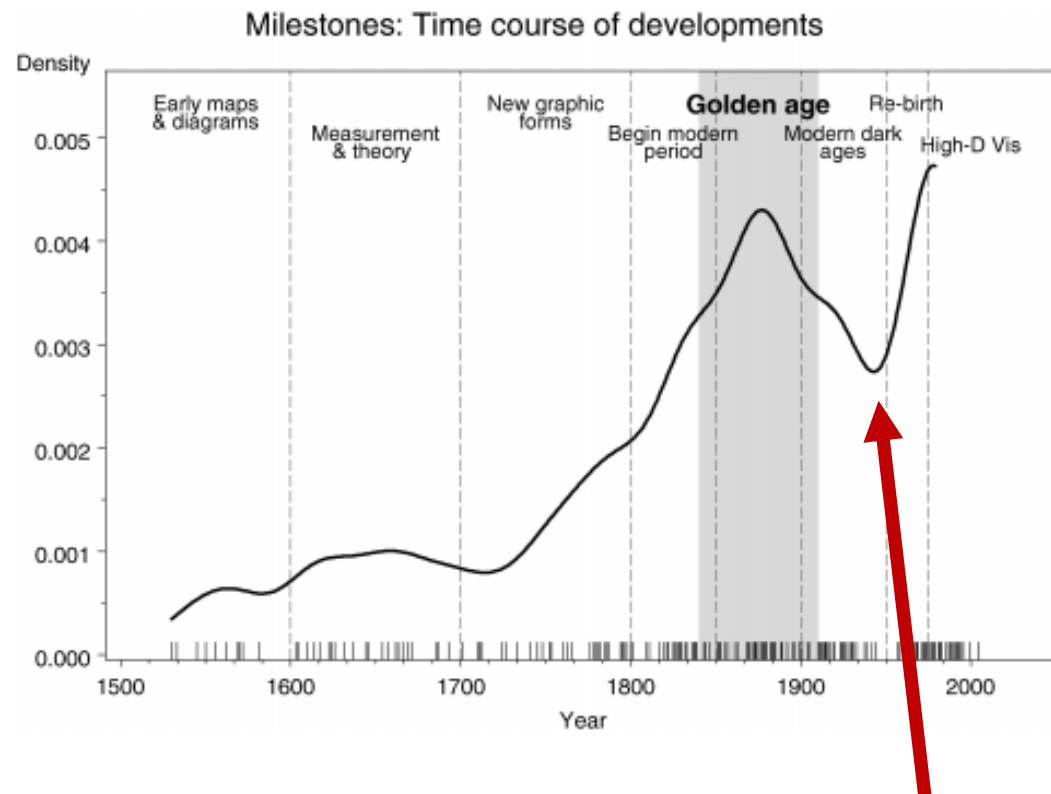# A very brief history of data visualization

Florence Nightingale (1820-1910)



Diagram of the causes of mortality in the army in the east

# A very brief history of data visualization

Francis Galton (1822-1911)



WEATHER CHART, MARCH 31, 1875.

The dotted lines indicate the gradations of barometric pressure. The variations of the temperature are marked by figures, the state of the sea and sky by descriptive words, and the direction of the wind by arrows—barbed and feathered according to its force. ⊙ denotes calm.

First weather map published in a newspaper (1875)

# A very brief history of data visualization

"Graphical dark ages" around 1950



Computer Age Statistical Inference, Efron and Hastie

# A very brief history of data visualization

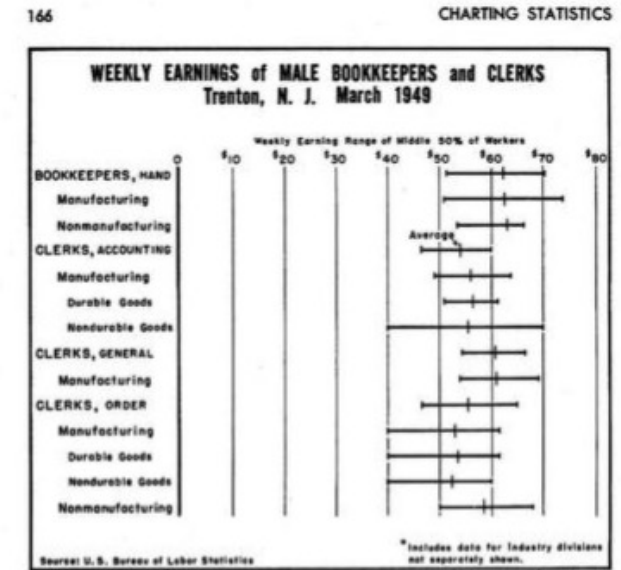Currently undergoing a "Graphical re-birth"

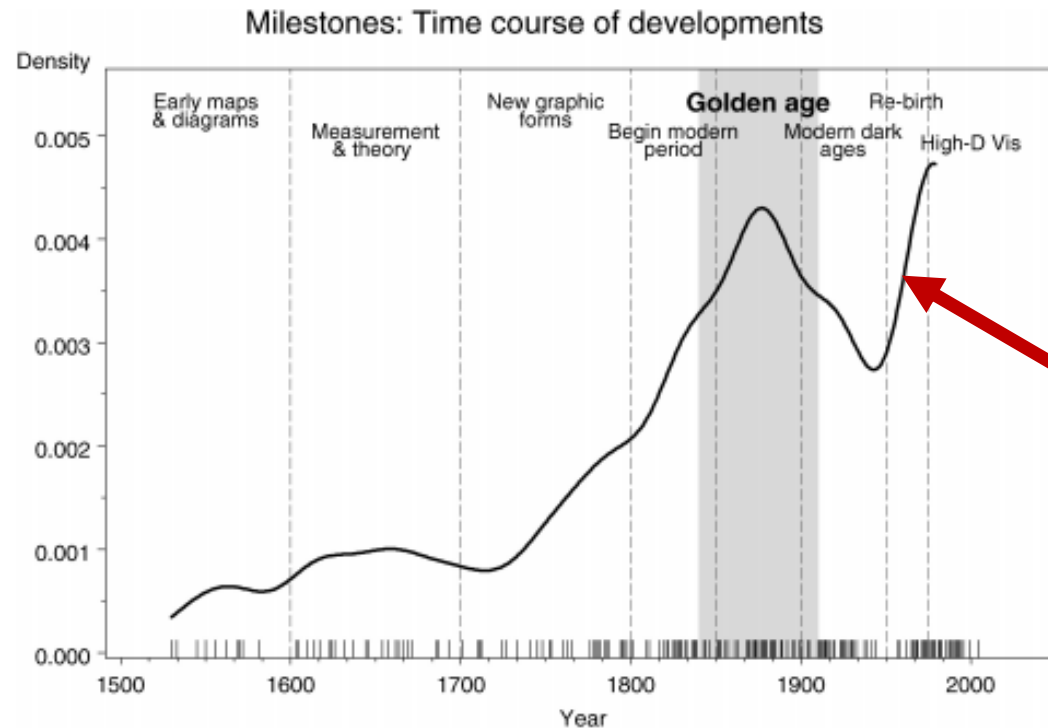Box plot



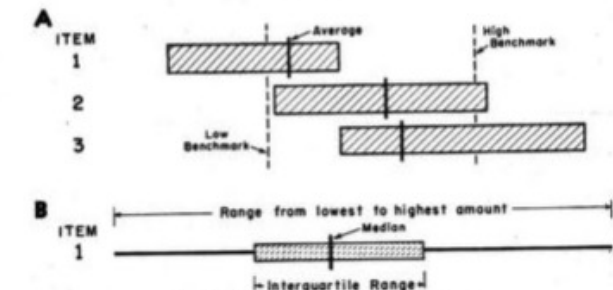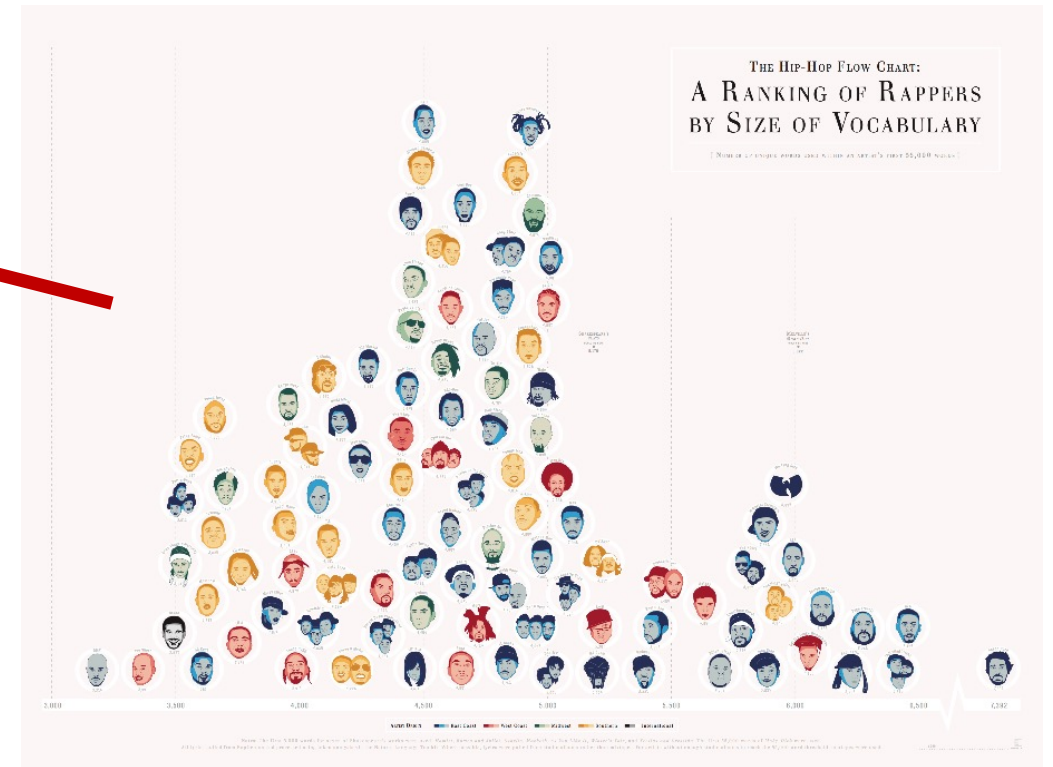Milestones: Time course of developments

Fig. 6-23. The range bar and symbol.
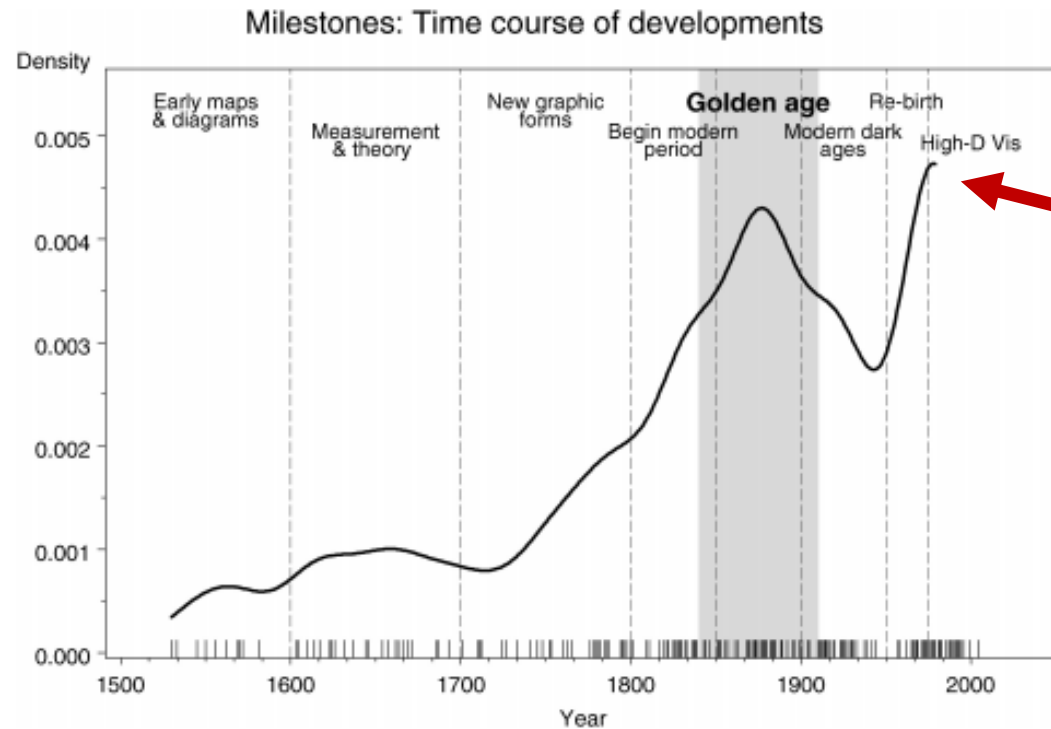
Fig. 6-24. Various uses of the range bar.
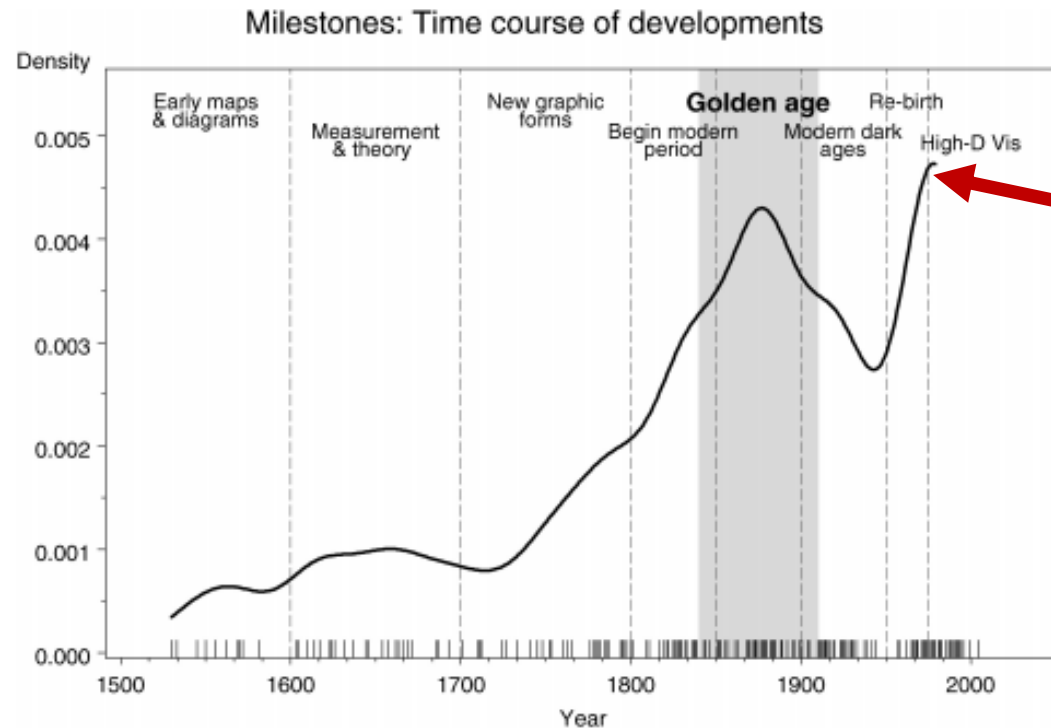
Spear 1952, Tukey 1970

# A very brief history of data visualization

Currently undergoing a "Graphical re-birth"

# A very brief history of data visualization

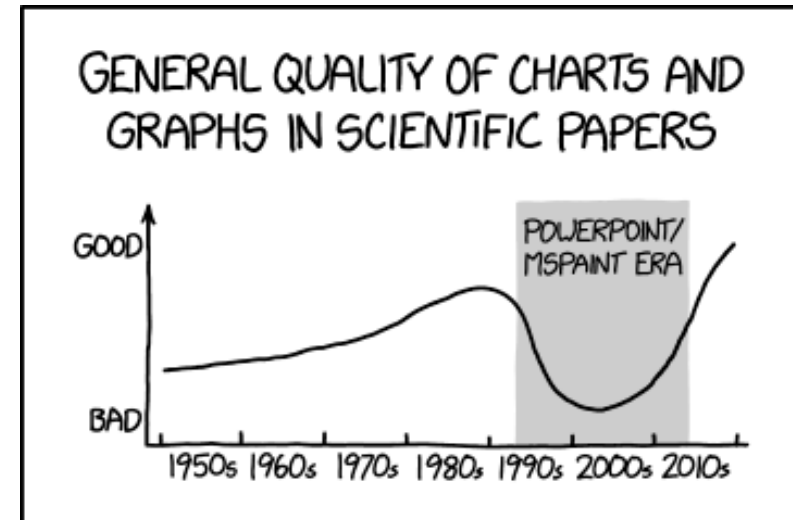Currently undergoing a "Graphical re-birth"



Hans Rosling's gapminder
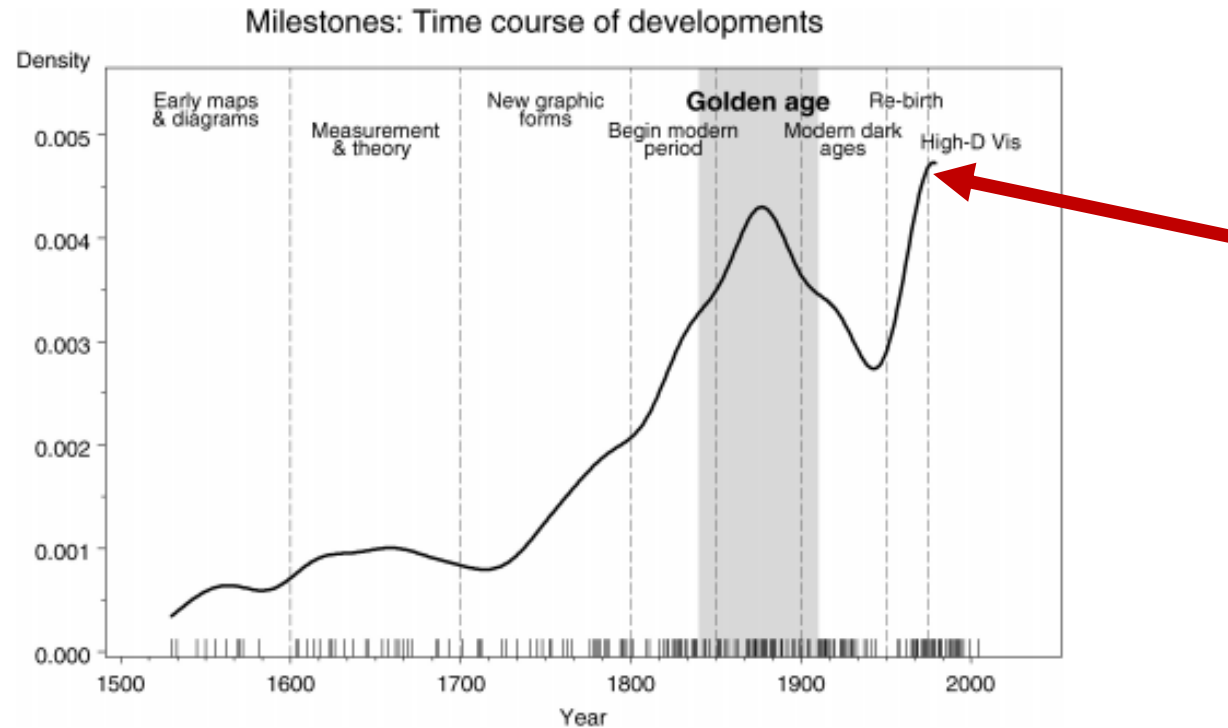
- [Simple version](#)
- [TV special effects](#)
- [Ted Talk](#)

# A very brief history of data visualization

Currently undergoing a "Graphical re-birth"
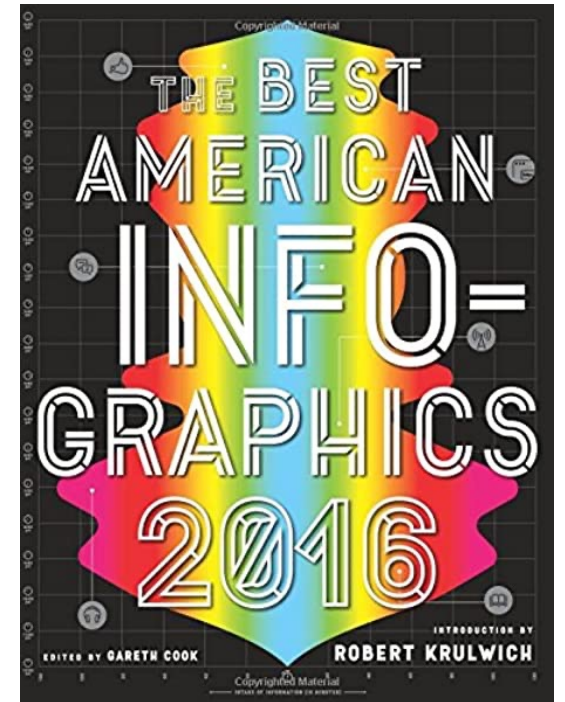
# Coming up on homework 5: find an interesting data visualization…

Homework 5 : Find an interesting data visualization
- https://www.reddit.com/r/dataisbeautiful/
- https://flowingdata.com/

We will do a little show and tell in class

# Visualizing data with matplotlib



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

- Matplotlib makes easy things easy and hard things possible.

Note: there are two different "interfaces" to matplotlib, that use slightly different syntax

- An explicit "Axes" interface
- An implicit "pyplot" interface
    - Just be aware if you are reading Stack overflow articles

Note: we will discuss seaborn soon...

We will use the pyplot interface (for now)

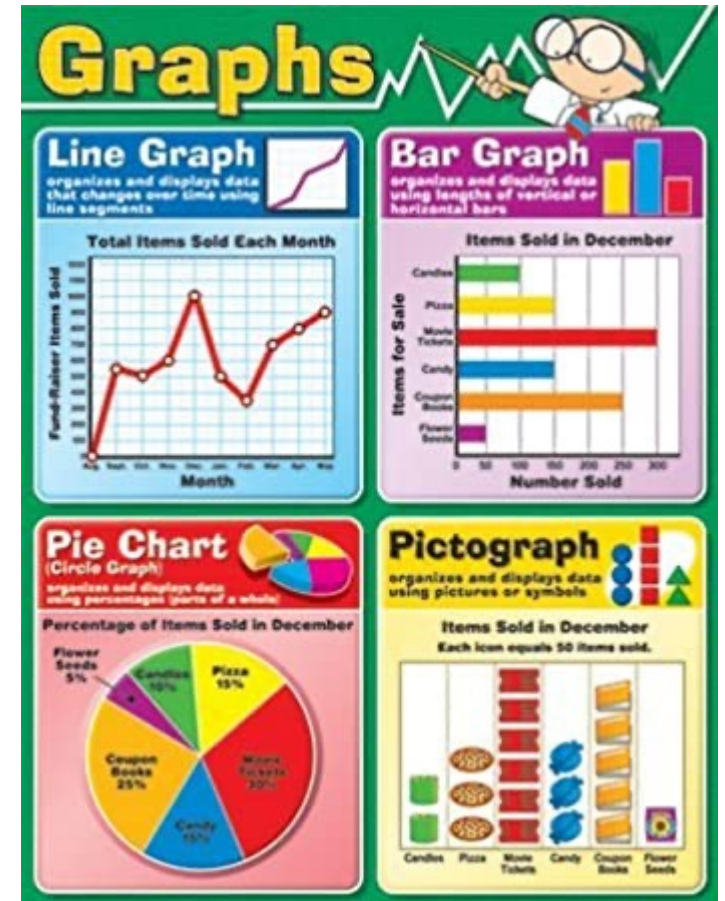    `import matplotlib.pyplot as plt`



seaborn

# Types of plots

There are many different ways to plot data, many of which I am sure you have seen before

The type of plot you choose will depend on the type of data you have and what you want to emphasize

- i.e., There are different types of plots of categorical data, a single quantitative variable, two quantitative variables, etc.
  - This will become even more apparent when we look at the seaborn package
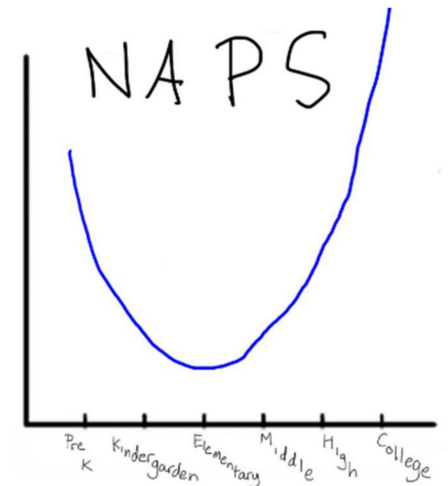
# Line graph
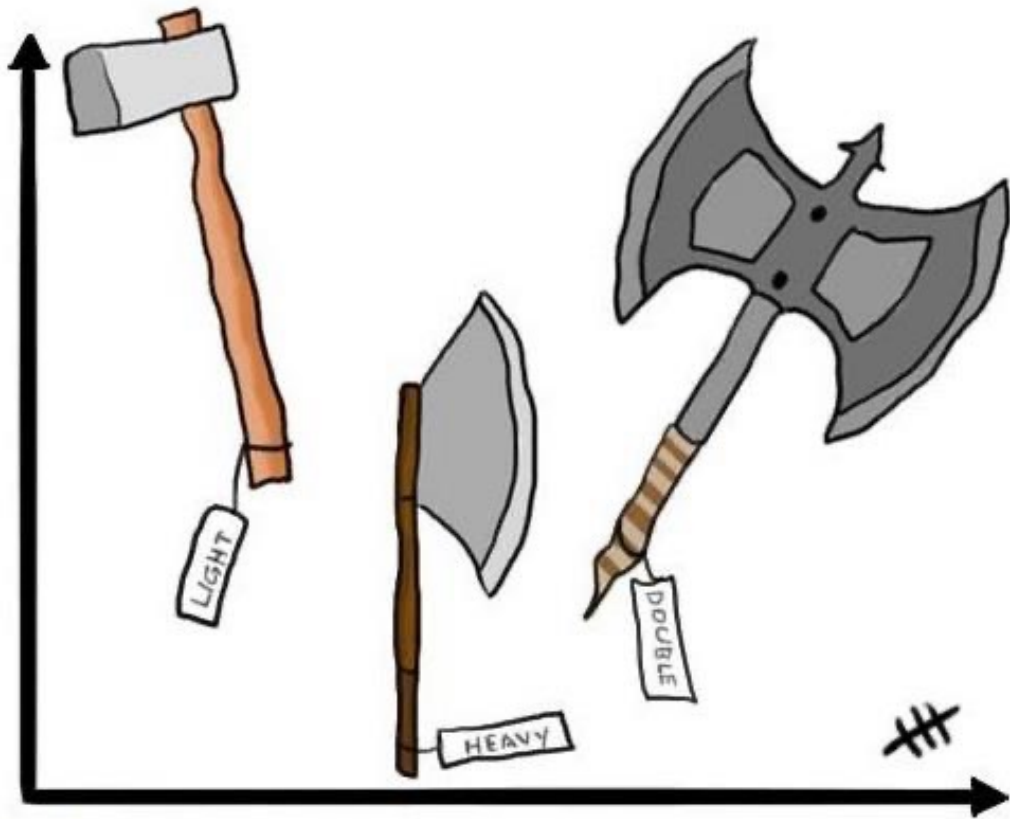
Line graph (line chart, or curve chart) displays information as a series of data points called "markers" connected by straight line segments.

We can create line graphs in matplotlib using:

- plt.plot(y)          # create a line graph with successive integers as the x-values
- plt.plot(x, y)         #  plots y as a function of x
- plt.plot(x, y, '-o')  # creates lines with circle markers

plt.ylabel("y label")

plt.xlabel("x label")

plt.title("my title")

plt.plot(x, y, label = "blah")

plt.legend()

Let's explore this in Jupyter!

# Visualizing quantitative data: histograms

Salaries of basketball players (in millions of dollars):
- 2.53,    18.6,    9.4,    21.7, …

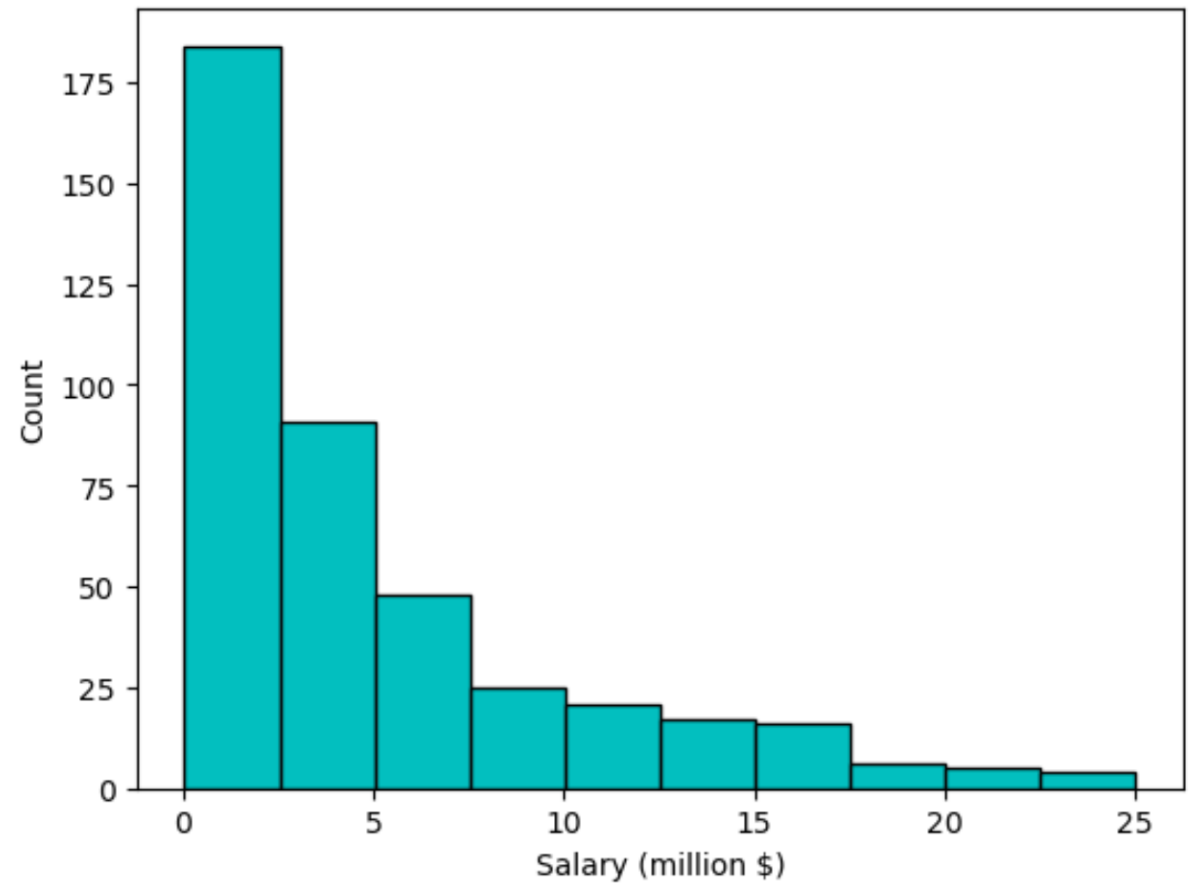To create a histogram we create a set of intervals
- 0-2.5,    2.5-5,    5-7.5,    …    20-22.5,    22.5-25.0

We count the number of points that fall in each interval

We create a bar chart where the height of the bars is the counts in each bin

# Histograms – countries life expectancy in 2007

| Life Expectancy | Frequency Count |
|---|---|
| (0 – 2.5] | 184 |
| (2.5 – 5] | 91 |
| (5 – 7.5] | 48 |
| (7.5 – 10] | 25 |
| (10 – 12.5] | 21 |
| (12.5 – 15] | 17 |
| (15 – 17.5] | 16 |
| (17.5 – 20] | 6 |
| (20 – 22.5] | 5 |
| (22.5 – 25] | 4 |



`plt.hist(data)`

Let's explore this in Jupyter!

# Five number summary

A **"Five Number Summary"** of quantitative data is defined as:

$$\text{(minimum, } Q_1\text{, median, } Q_3\text{, maximum)}$$

- $Q_1$ is the 25th percentile - i.e., the value such that 25% of the data is less than this value
- $Q_3$ is the 75th percentile - i.e., the value such that 75% of the data is less than this value

The Five number summary roughly splits the data into 4 equal parts

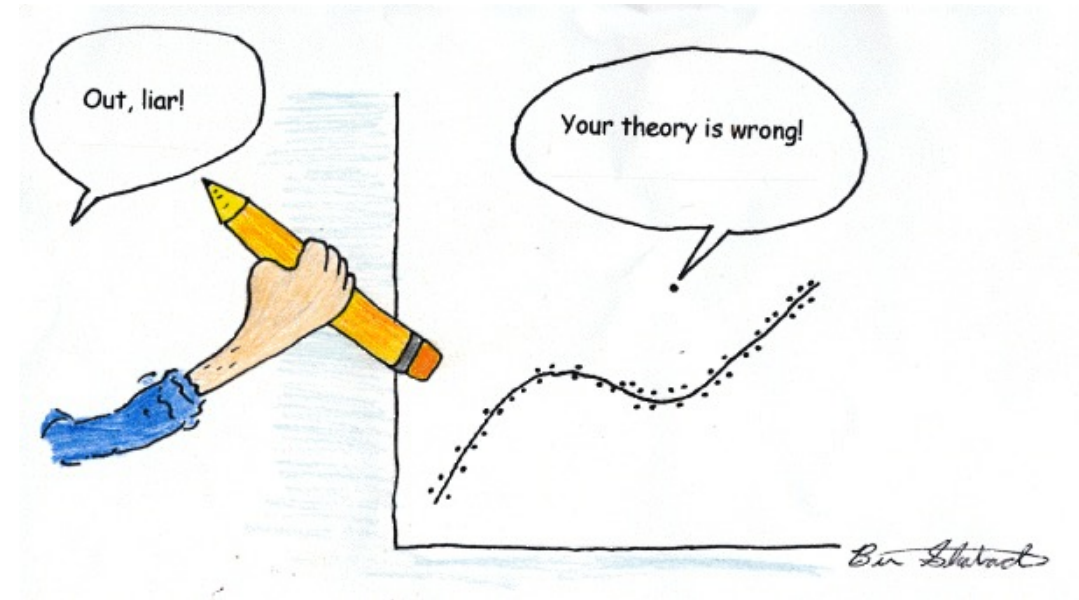The **Interquartile range (IQR)** is defined as $Q_3 - Q_1$

# Detecting of outliers

As a rule of thumb, we call a data value
an **outlier** if it is:

Smaller than: $Q_1 - 1.5 * IQR$

Larger than: $Q_3 + 1.5 * IQR$

# Boxplots

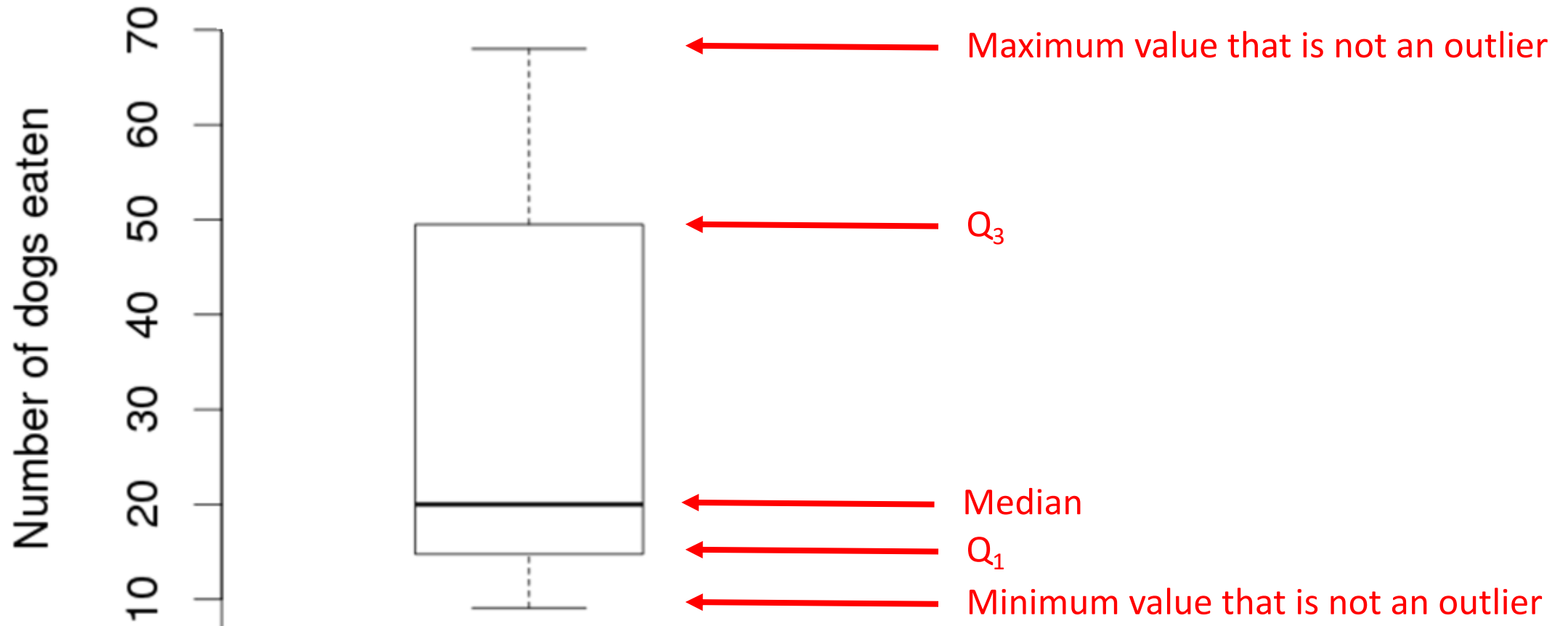A **boxplot** is a graphical display of the 5 number summary and consists of:

1. Drawing a box from $Q_1$ to $Q_3$

2. Dividing the box with a line (or dot) drawn at the median

3. Draw a line from each quartile to the most extreme data value that is not and outlier

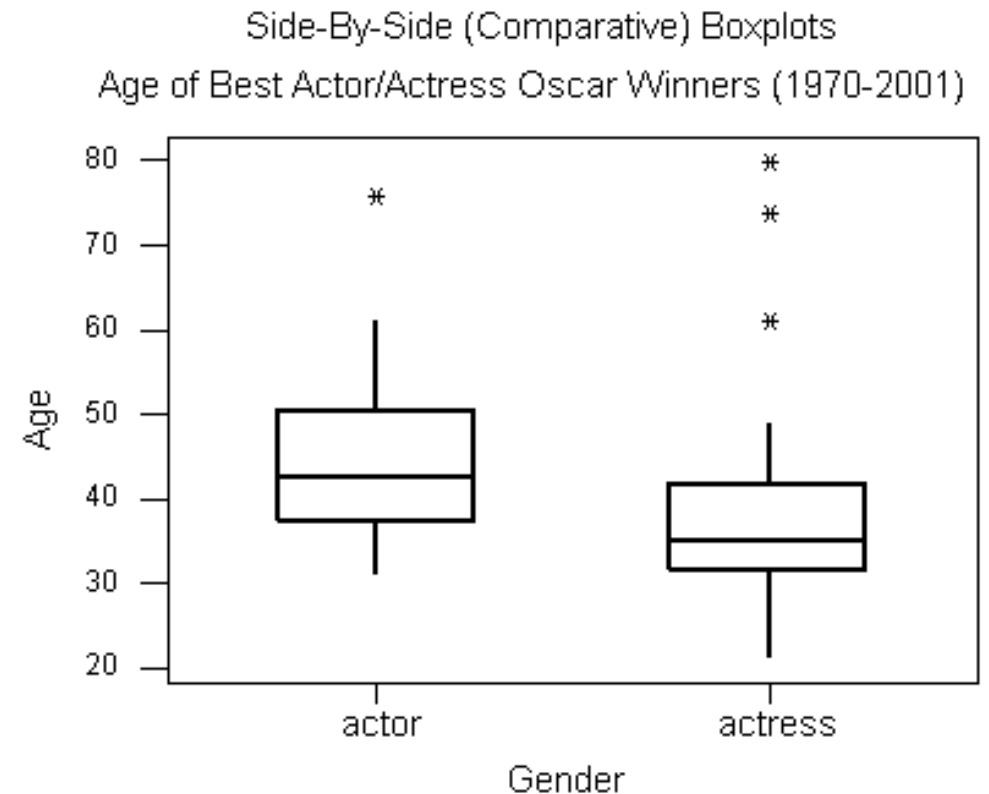4. Draw a dot/asterisk for each outlier data point.

# Box plot of the number of hot dogs eaten by the men's contest winners 1980 to 2010

# Comparing quantitative variables across categories

Often one wants to compare quantitative variables across categories

**Side-by-Side** graphs are a way to visually compare quantitative variables across different categories.



Side-By-Side (Comparative) Boxplots
Age of Best Actor/Actress Oscar Winners (1970-2001)

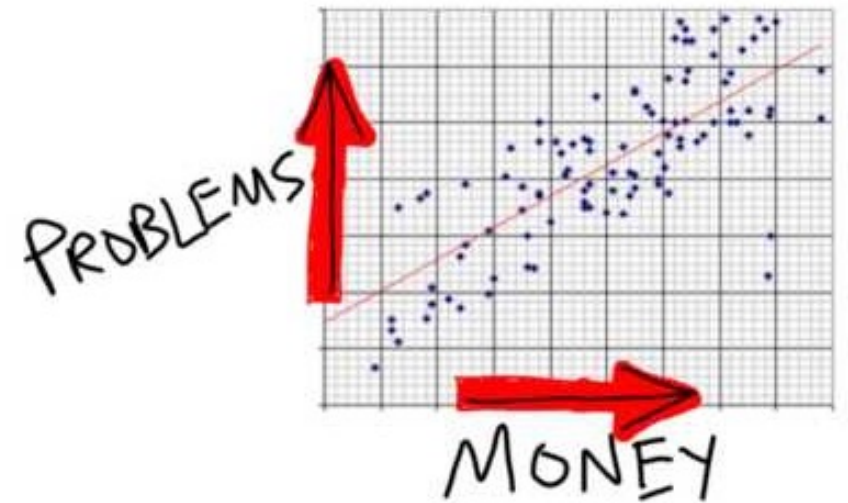Let's explore this in Jupyter!

`plt.boxplot(data)`

# Scatterplots

A **scatterplot** graphs the relationship between two variables

> Each axis represents the value of one variables

> Each point the plot shows the value for the two variables for a single data case

If there is an explanatory and response variable, then the explanatory variable is put on the x-axis and the response variable is put on the y-axis.



PROBLEMS

MONEY

# Scatterplots

There are two ways to create scatter plots in matplotlib:

1. Using plt.plot(x, y, '.')

2. Using plot.scatter(x, y)
   - This function has additional useful arguments such as:
     - s:  specified the size of each point
     - color:  specifies the color of each point
     - marker: specifies the shape of each point

Let's explore this in Jupyter!

# Bar plots and pie charts

Bar plots and pie charts are used to plot **categorical data**

Bar plots, the heights of the bars indicate the number of items in each category

In pie charts, the angle of each segment indicates the proportion of items in each category

Let's explore this in Jupyter!



**REAL Bar Chart**

% of Drinks Ordered

- 22.7% Beer
- 17.1% Wine
- 9.9% Martini
- 8.7% Margarita
- 7.2% Ice Tea
- 6.1% Rum & Coke
- 6.0% Mai Tai
- 4.6% Gin & Tonic

**World's Most Accurate Pie Chart**

- Pie I have eaten
- Pie I have not yet eaten

# Subplots

Matplotlib makes it easy to create multiple subplots within a larger figure

2 columns

1 row

plt.subplot(1, 2, 1);  ← plot on the first subplot
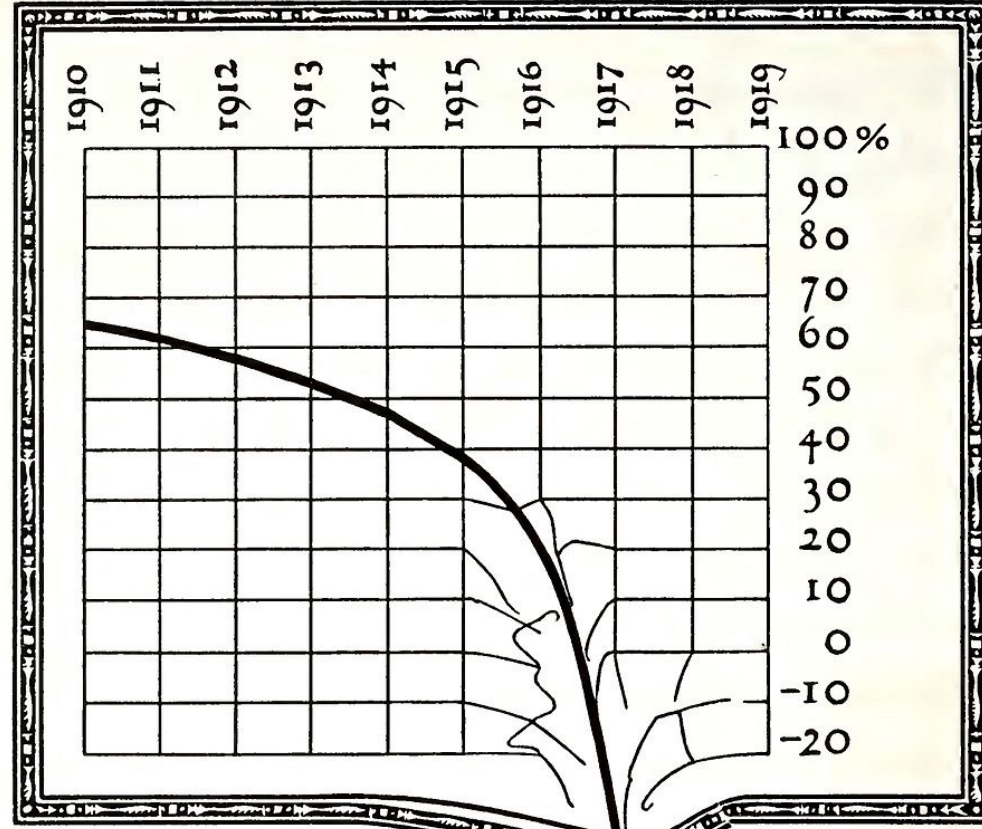plt.plot(x1, y1);

plt.subplot(1, 2, 2);  ← plot on the second subplot
plt.plot(x2, y2);

Let's explore this in Jupyter!

A chart showing the percentage of excellence in the physical properties of books published since 1910.