

# INTRODUCTION



To DATA SCIENCE

# Introduction to Data Science

## Collecting, Exploring and Cleaning Data - Class 3

Giora Simchoni

[gsimchoni@gmail.com](mailto:gsimchoni@gmail.com) and add #intro2ds in subject

Stat. and OR Department, TAU

INTRODUCTION



To DATA SCIENCE

# Common Data Formats in Data Science

INTRODUCTION



TO DATA SCIENCE

# CSV: Comma Separated Values

```

1 id,drugName,condition,review,rating
2 206461,Valsartan,Left Ventricular Dysfunction,"""It has n
3 95260,Guanfacine,ADHD,"""My son is halfway through his fo
4 We have tried many different medications and so far this
5 92703,Lybrel,Birth Control,"""I used to take another oral
6 The positive side is that I didn't have any other si
7 138000,Ortho Evra,Birth Control,"""This is my first time
8 35696,Buprenorphine / naloxone,Opiate Dependence,"""Subox

```

```

import pandas as pd

df = pd.read_csv('..../datasets/drugs.csv')

print(df.head())

```

	id	drugName	condition	review	rating
0	206461	Valsartan	Left Ventricular Dysfunction	"It has no side effect, I take it in combinati...	9
1	95260	Guanfacine	ADHD	"My son is halfway through his fourth week of ...	8
2	92703	Lybrel	Birth Control	"I used to take another oral contraceptive, wh...	5
3	138000	Ortho Evra	Birth Control	"This is my first time using any form of birth...	8
4	35696	Buprenorphine / naloxone	Opiate Dependence	"Suboxone has completely turned my life around...	9

# JSON: JavaScript Object Notation

```
1 | {"firstname": "John",
2 |   "lastname": "Smith",
3 |   "age": "27",
4 |   "address": {"streetAddress": "21 2nd Street",
5 |     "city": "New York",
6 |     "state": "NY",
7 |     "postalCode": "10021-3100"},
8 |     "children": [],
9 |     "spouse": null}
```

```
import json

data = dict()

with open('../datasets/test.json') as f:
    data=json.load(f)

data
```

```
{'firstname': 'John',
'lastname': 'Smith',
'age': '27',
'address': {'streetAddress': '21 2nd Street',
'city': 'New York',
'state': 'NY',
'postalCode': '10021-3100'},
'children': [],
'spouse': None}
```

# Plain Text

```
If you should go skating  
On the thin ice of modern life  
Dragging behind you the silent reproach  
Of a million tear-stained eyes  
Don't be surprised when a crack in the ice  
Appears under your feet.
```

```
with open('../datasets/test.txt') as f:  
    lines = f.readlines(1000)
```

```
lines
```

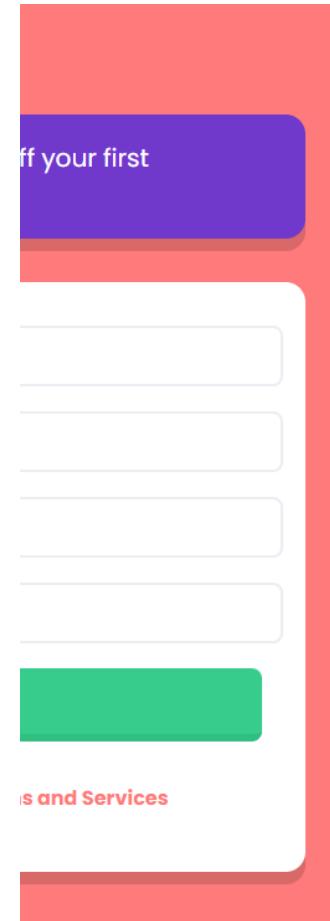
```
['If you should go skating,\n',  
'On the thin ice of modern life,\n',  
'Dragging behind you the silent reproach,\n',  
'Of a million tear-stained eyes,\n',  
 "Don't be surprised when a crack in the ice,\n",  
 'Appears under your feet.]
```

# HTML

```

190 <body>
191   <main>
192     <!-- intro section -->
193     <section class="intro">
194       <h1 class="title">Black Friday Deals</h1>
195       <p>Get up to 50% off on all our products and services. Hurry up, the offer ends in 24 hours.</p>
196     </section>
197
198     <!-- sign-up section -->
199     <section class="sign-up">
200       <p class="sign-up-para">Sign up for our newsletter and get 10% off your first purchase</p>
201       <!-- the form itself -->
202       <form class="sign-up-form">
203         <div class="form-input">
204           <input type="text" name="first-name" id="first-name" placeholder="First Name" required>
205           <span>!</span>
206           <p class="warning">First name cannot be empty</p>
207         </div>
208
209         <div class="form-input">
210           <input type="text" name="last-name" id="last-name" placeholder="Last Name" required>
211           <span>!</span>
212           <p class="warning">Last name cannot be empty</p>
213         </div>
214
215         <div class="form-input">
216           <input type="email" name="email" id="email" placeholder="Email Address" required>
217           <span>!</span>
218           <p class="warning">Looks like this is not an email</p>
219         </div>
220
221         <div class="form-input">
222           <input type="password" name="Password" id="Password" placeholder="Password" required>
223           <span>!</span>
224           <p class="warning">Password cannot be empty</p>
225         </div>
226
227         <input class="submit-btn" type="submit" value="Claim your offer">
228         <p class="form-term">By clicking the button, you are agreeing to our <span>Terms and Services</span> </p>
229     </form>

```



# HTML

```
from bs4 import BeautifulSoup

with open('../datasets/test.html') as f:
    soup = BeautifulSoup(f, 'html.parser')

print(soup.prettify())

<!DOCTYPE html>
<html lang="en">
    <head>
        <meta charset="utf-8"/>
        <meta content="IE=edge" http-equiv="X-UA-Compatible"/>
        <meta content="width=device-width, initial-scale=1.0" name="viewport"/>
    <title>
        Sample HTML Code - NewsLetter Form
    </title>
    <style>
        @import url("https://fonts.googleapis.com/css2?family=Poppins:wght@400;500;600;700&display=swap");

        body {
            display: flex;
            justify-content: center;
            padding: 3rem 0;
            font-family: "Poppins", sans-serif;
            font-size: 1rem;
            color: white;
            background-color: #ff7a7a;
        }
    </style>

```

# Collecting Data

INTRODUCTION



To DATA SCIENCE

# Where do(es) data come from?

- Then: Manual Curation
- Now: Automatic Curation

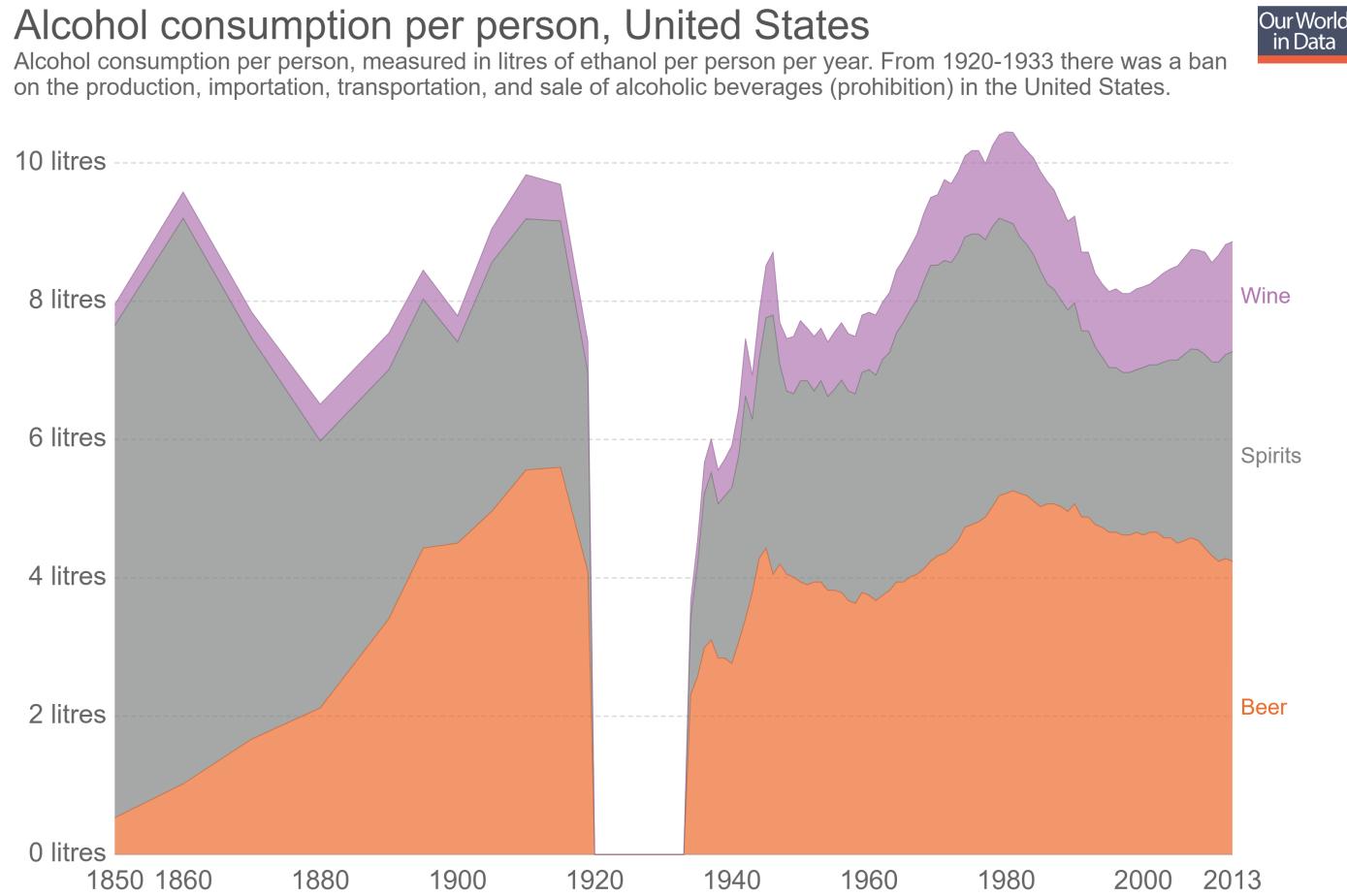
# Then: Manual Curation (I)

If your parents have not taken note, anywhere, of how tall you were at the age of 1 - we may never be able to extract this information.



# Then: Manual Curation (II)

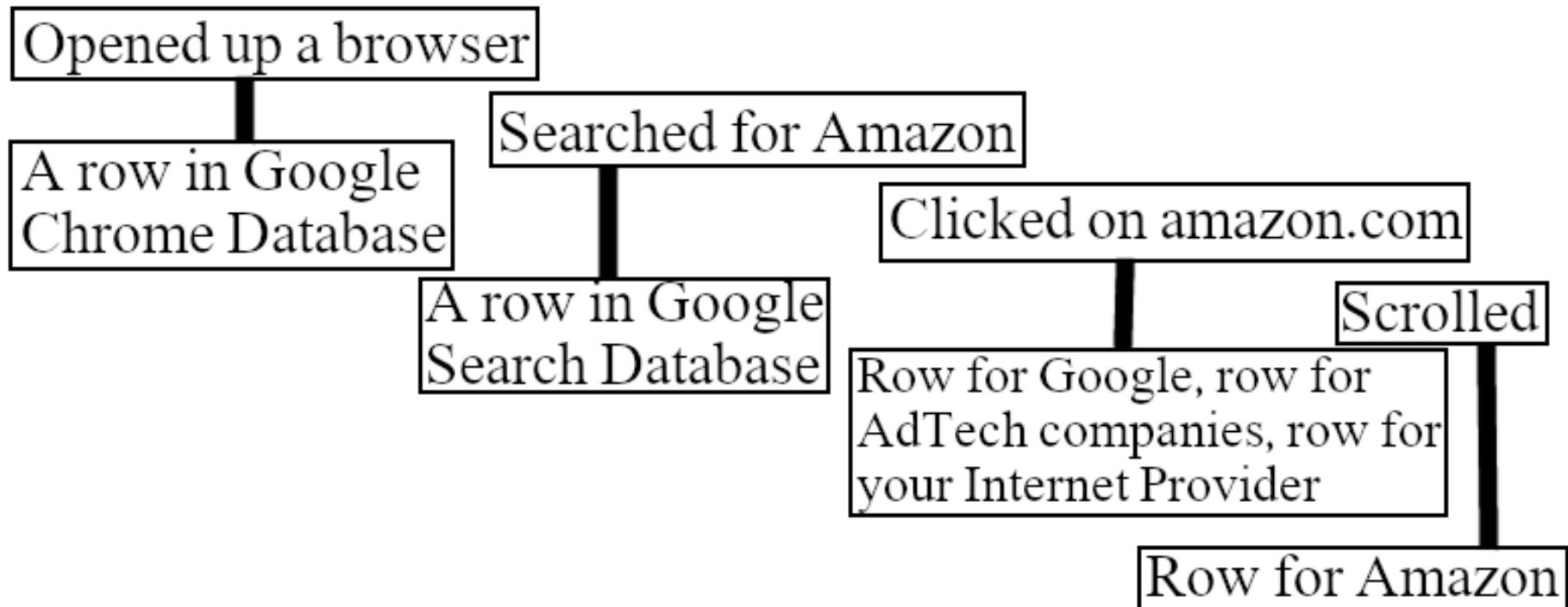
If the US government had not seen fit to estimate and record the level of alcohol consumption of its citizens, we would never have known.



# Now: Automatic Curation (I)

Have you ever opened up an Internet browser, searched for “Amazon”, clicked on [amazon.com](http://amazon.com) and scrolled around to check the price of a T-shirt? You don’t have to be logged in. You don’t have to buy. You are data.

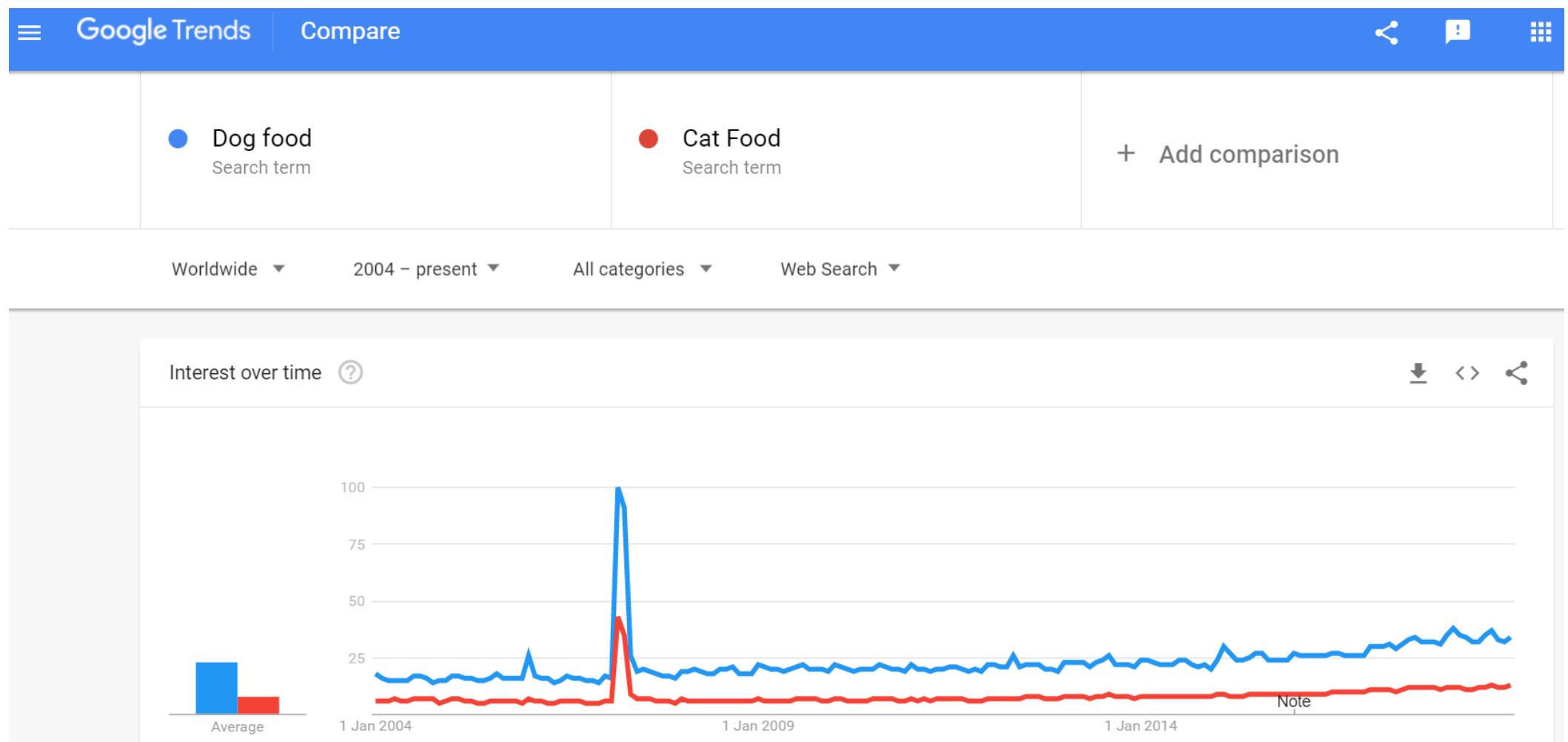
You are Data.



# Web Scraping

- Public APIs
- Beautiful Soup

# Public APIs (I)



# Beautiful Soup

You may not even need an API. Would you look to have all the [Beatles](#) records in a single table on your machine?

The following code scrapes the Wikipedia page for the and creates a table, out of “thin air”. See more advanced examples in recitation.

```
1 from bs4 import BeautifulSoup
2 import requests
3 import re
4 import pandas as pd
5
6 url = 'https://en.wikipedia.org/wiki/The_Beatles_discography'
7 r = requests.get(url)
8 soup = BeautifulSoup(r.content, 'html.parser')
```



```
1 albums = dict()
2 id = 0
3 albums[id] = dict()
4 tables = soup.find_all('table')
5 for table in tables:
6     caption = table.find('caption')
7     if caption is not None:
8         header = caption.get_text()
9         if re.match(re.compile('^List of(.+?)albums'), header):
10             rows = table.find_all('tr')
11             for row in rows:
12                 title_col = row.find('th')
13                 if title_col is not None and 'scope' in title_col.attrs and\
14                     title_col.attrs['scope'] == 'row':
15                     title_cell = title_col.find('a')
16                     if title_cell is not None and title_cell.attrs is not None and\
17                         'title' in title_cell.attrs:
18                         albums[id]['name'] = title_cell.attrs['title']
19                         release_col = row.find('td')
20                         release_date, release_label = get_release_details(release_c
```

```
albums_df = pd.DataFrame.from_dict(albums, orient ='index')
albums_df.head(5)
```

	<b>name</b>	<b>release_date</b>	<b>release_label</b>
0	Please Please Me	22 March 1963	Parlophone
1	With the Beatles	22 November 1963	Parlophone (UK), Capitol (Canada), Odeon (France)
2	A Hard Day's Night (album)	10 July 1964	Parlophone
3	Beatles for Sale	4 December 1964	Parlophone
4	Help!	6 August 1965	Parlophone

# Small Data, Big Data

INTRODUCTION



To DATA SCIENCE

# What's in a name?

These definitions are constantly changing.

- “Everything processed in Excel is small data.” ([Rufus Pollock, The Guardian](#))
- “[Big Data] is data so large it does not fit in main memory” (Leskovec et al., Mining of Massive Datasets)

Or maybe we should define the size of our data according how easy it is to process and understand it?

- “[Small Data is] data that has small enough size for human comprehension.” ([jWork.ORG](#))
- “data sets that are too large or complex for traditional data-processing application software to adequately deal with” ([Wikipedia](#))

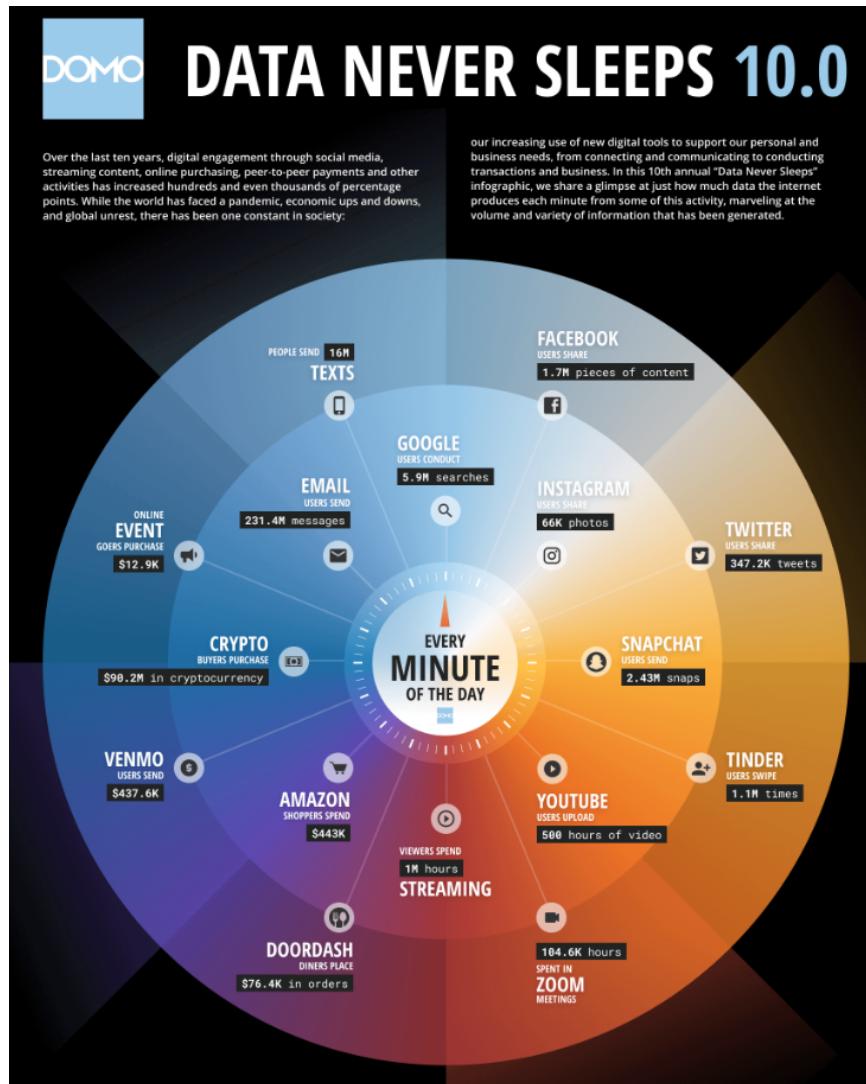
# What's in a name? (II)

The actual definition should probably merge both of the above.

- Excel can fit 1M rows, 16K columns of double numbers. Try loading a matrix such as this into Matlab, Python or R, multiply and invert it - you can't. So isn't that Big?
- Facebook generates 4 Petabytes of data, daily. That's 4K Terabytes or 4M Gigabytes. ([Brandwatch.com](#)) But a Facebook Data Scientist in daily life typically needs only a copy of some of these data, which fits in her PC. Isn't that small?

# Web data is Big Data

We can all agree *this* is big: ([Domo.com](https://domo.com/data-never-sleeps))



# Exploring Data: Basic Plots

INTRODUCTION

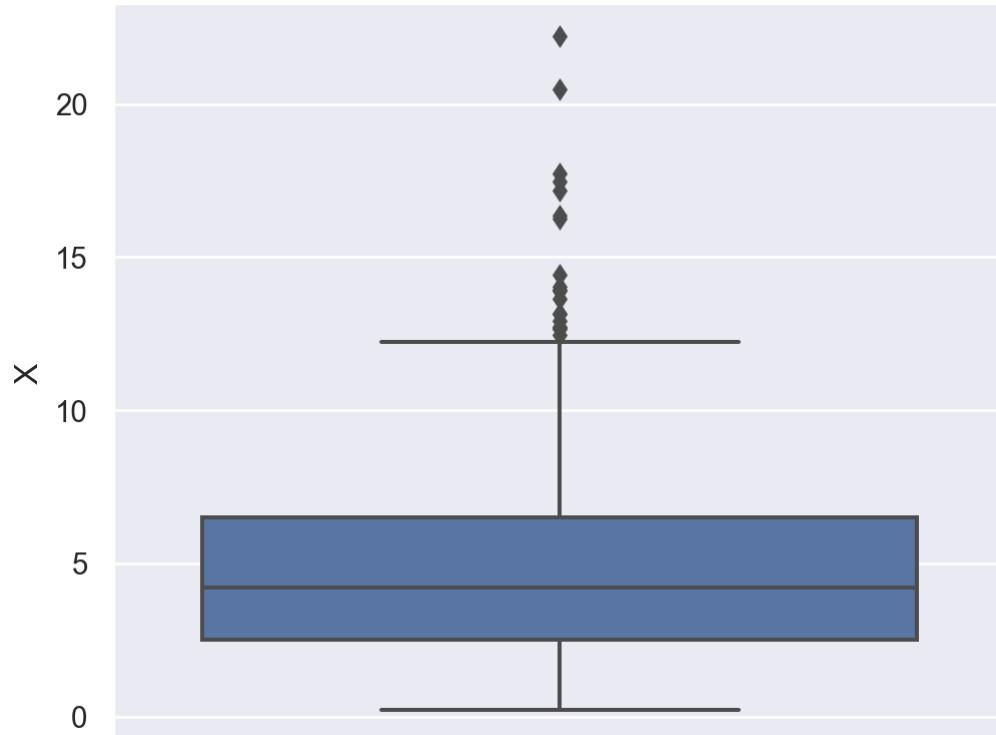


To DATA SCIENCE

# Boxplot

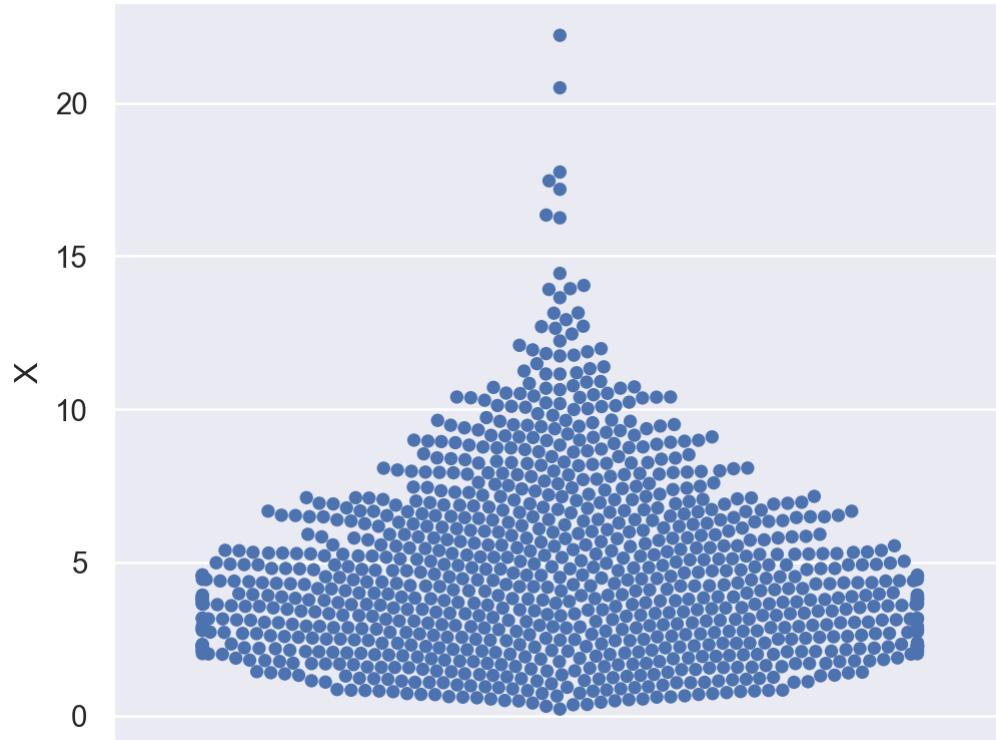
```
import numpy as np  
import seaborn as sns  
import matplotlib.pyplot as plt
```

```
sns.boxplot(y = X)  
plt.ylabel('X')  
plt.show()
```



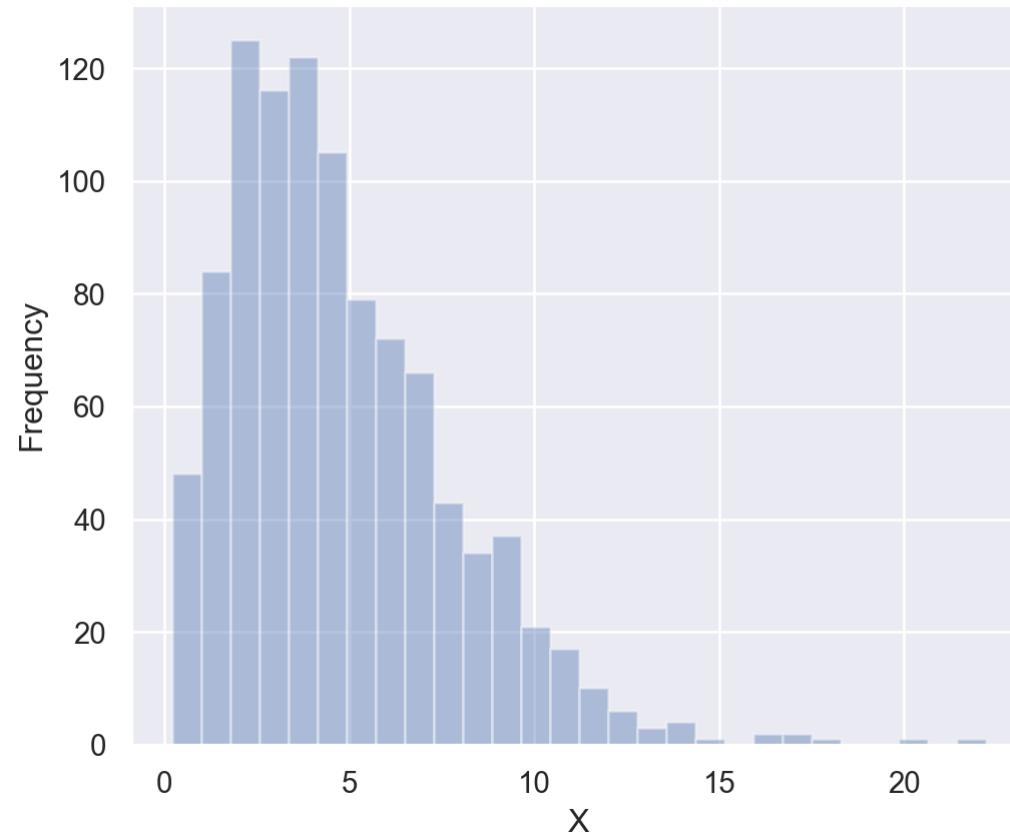
# Swarmplot

```
sns.swarmplot(y = X)  
plt.ylabel('X')  
plt.show()
```



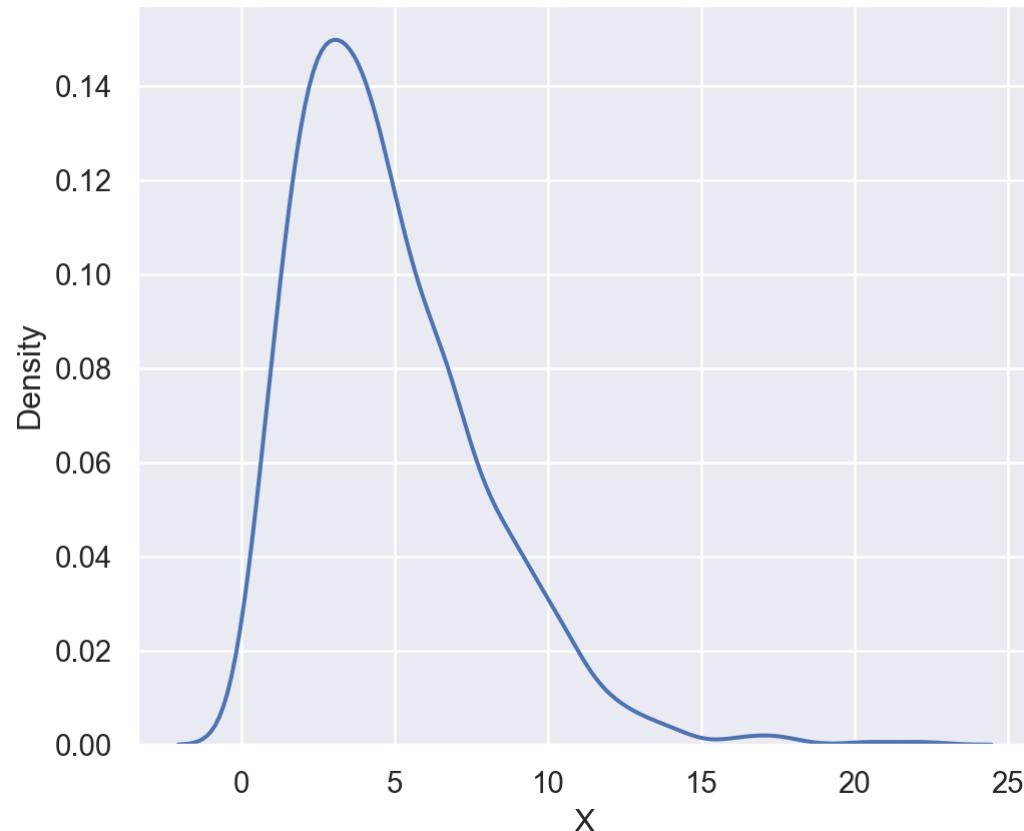
# Histogram

```
sns.distplot(X, kde = False)  
plt.ylabel('Frequency')  
plt.xlabel('X')  
plt.show()
```



# Density Plot

```
sns.distplot(X, hist = False)  
plt.ylabel('Density')  
plt.xlabel('X')  
plt.show()
```



# Density plot: kernel density estimation / convolution

The way to get a smooth estimate of the distribution in density plot is by defining a kernel which “smoothes” the data. Mathematically we define a kernel weight function  $w : \mathbb{R} \rightarrow \mathbb{R}^+$  as:

1. Non-negative and symmetric:  $w(x) = w(-x)$
2. Integrates to 1:  $\int_{\mathbb{R}} w(x)dx = 1$

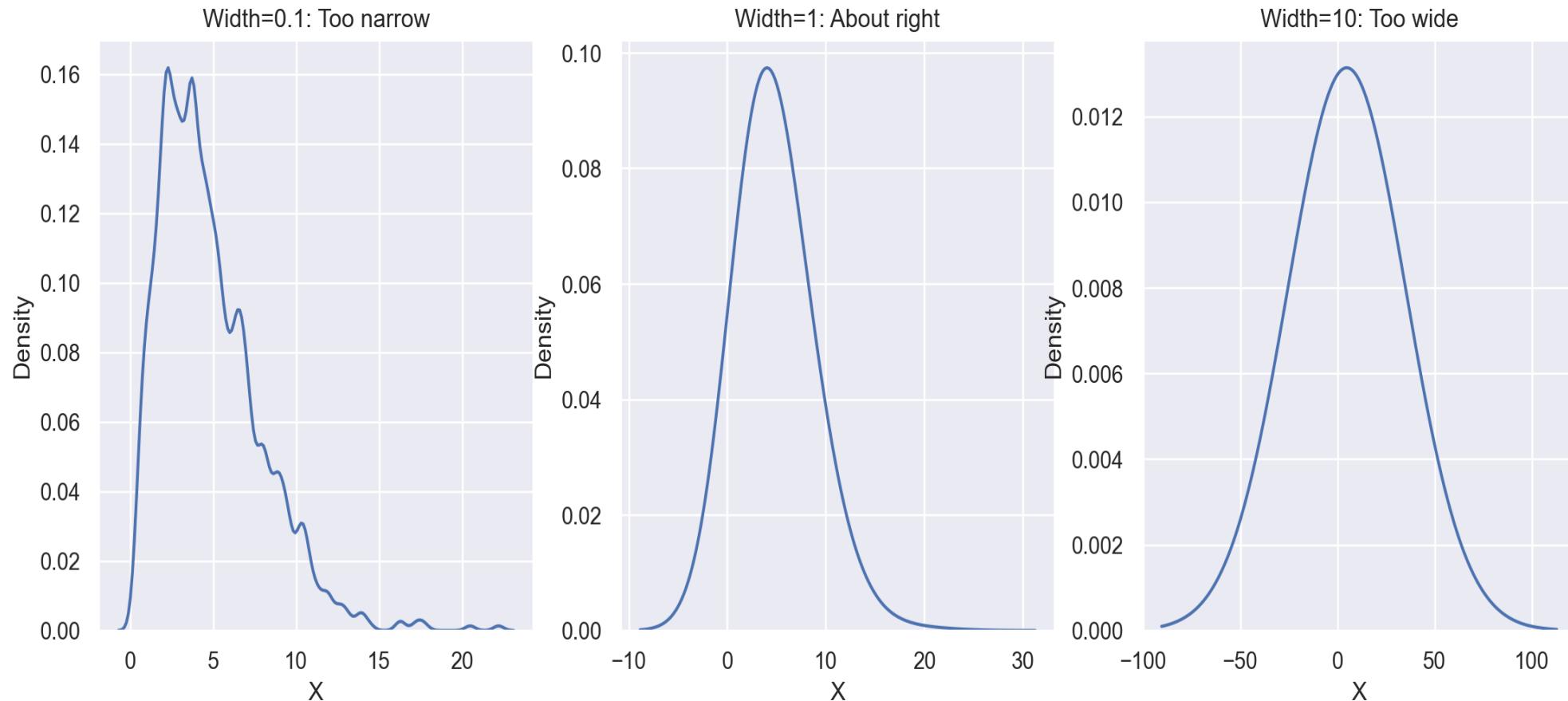
And then the density kernel estimate is:  $J(x) = \frac{1}{n} \cdot \sum_{i=1}^n w(x_i - x)$ .

Nice property:  $\int_{\mathbb{R}} J(x)dx = 1$

- Wide  $w$ : smooth estimate, but it may not reflect the real data
- Narrow  $w$ : very non-smooth description

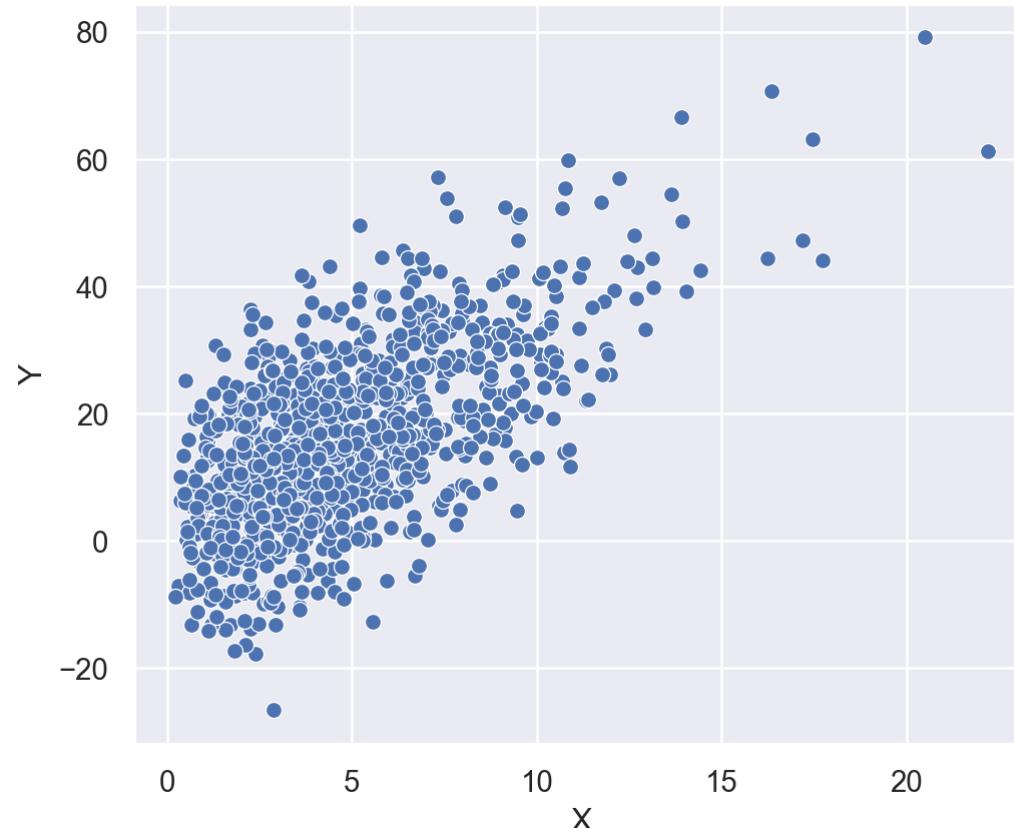
For example smoothing this same dataset with too narrow or too wide window:

► Code



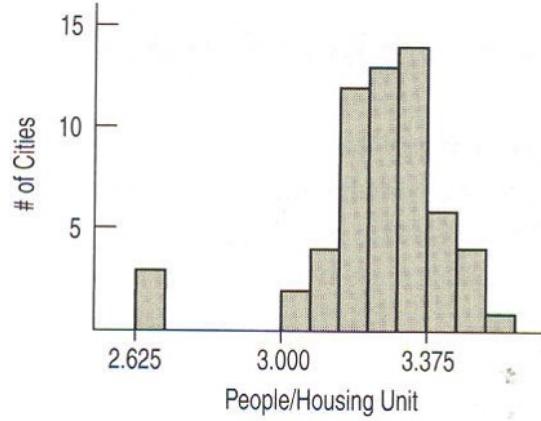
# Scatterplot

```
Y = X * 3 + 2 + np.random.normal(0, 10, 1000)
sns.scatterplot(X, Y)
plt.ylabel('Y')
plt.xlabel('X')
plt.show()
```

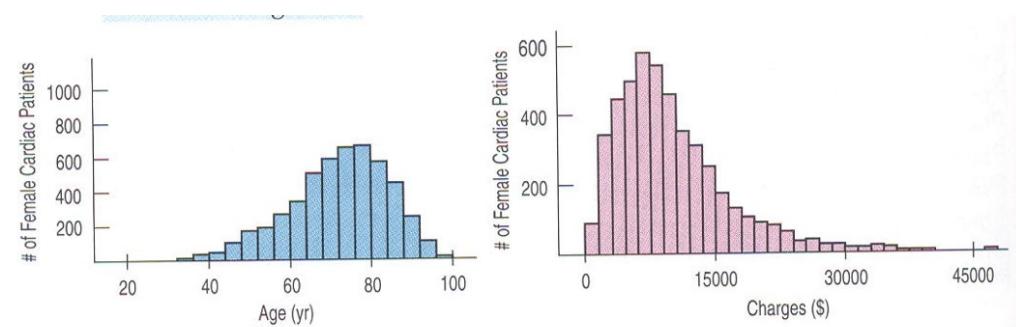


# What can we learn from simple plots?

Look at outliers:



See the shape and tail direction:



Age at heart attack (left) and cost of hospitalization (right)

# Exploring Data: Summary Statistics

INTRODUCTION



To DATA SCIENCE

# Location

“Where is this X located? Where is the central mass?”

- Mean of empirical distribution (=average):

$$\text{Mean}(X) = \bar{X} = \frac{1}{N} \sum_{i=1}^N X_i$$

- Median:

$$\text{Med}(X) = m \text{ s.t. } P(X \leq m) = P(X \geq m) = 0.5$$

- Mode:

$$\text{Mode}(X) = \text{Most frequent value in a dataset}$$

```
mean = np.mean(X)
median = np.median(X)
hist, _ = np.histogram(X, bins=range(20))
mode = list(range(20))[hist.argsort()[-1][0]]
plt.figure(figsize= (4,4))
sns.distplot(X, bins = range(20), kde = False)
plt.plot([mean, mean], [0, 160], linewidth=2, color='r')
plt.plot([median, median], [0, 160], linewidth=2, color='g')
plt.plot([mode, mode], [0, 160], linewidth=2, color='b')
plt.legend({'Mean':mean,'Median':median,'Mode':mode})
plt.xlabel('X')
plt.ylabel('Frequency')
plt.show()
```

# Dispersion

“Is X widely spread out? Does it concentrate narrowly around the mean?”

- Quantiles/Percentiles:

$$Q(X, q) = v \text{ s.t. } P(X \leq v) = 1 - P(X \geq v) = q$$

- Range:

$$\text{Range}(X) = \text{Max}(X) - \text{Min}(X)$$

- Inter-Quartile-Range:

$$IQR(X) = Q(X, 0.75) - Q(X, 0.25)$$

# Dispersion

- (Empirical) Variance:

$$Var(X) = \frac{1}{N} \sum_{i=1}^N (X_i - Mean(X))^2$$

- Standard Deviation:

$$STD(X) = \sqrt{Var(X)}$$

```
print(f'90th percentile: {np.percentile(X, 90) :.2f}')
print(f'Range: {np.max(X) - np.min(X) :.2f}')
print(f'IQR: {np.percentile(X, 75) - np.percentile(X, 25) :.2f}')
print(f'Variance: {np.var(X) :.2f}')
print(f'Standard Deviation: {np.std(X) :.2f}')
```

90th percentile: 8.95

Range: 22.00

IQR: 3.97

Variance: 9.12

Standard Deviation: 3.02

# Shape

“Is X symmetric or not? How ‘tailed’ is it?”

- Skewness:

$$Skew(X) = \frac{1}{N} \frac{\sum_{i=1}^N (X_i - Mean(X))^3}{STD(X)^3}$$

```
from scipy import stats  
  
print(f'Skewness: {stats.skew(X) :.2f}')
```

Skewness: 1.22

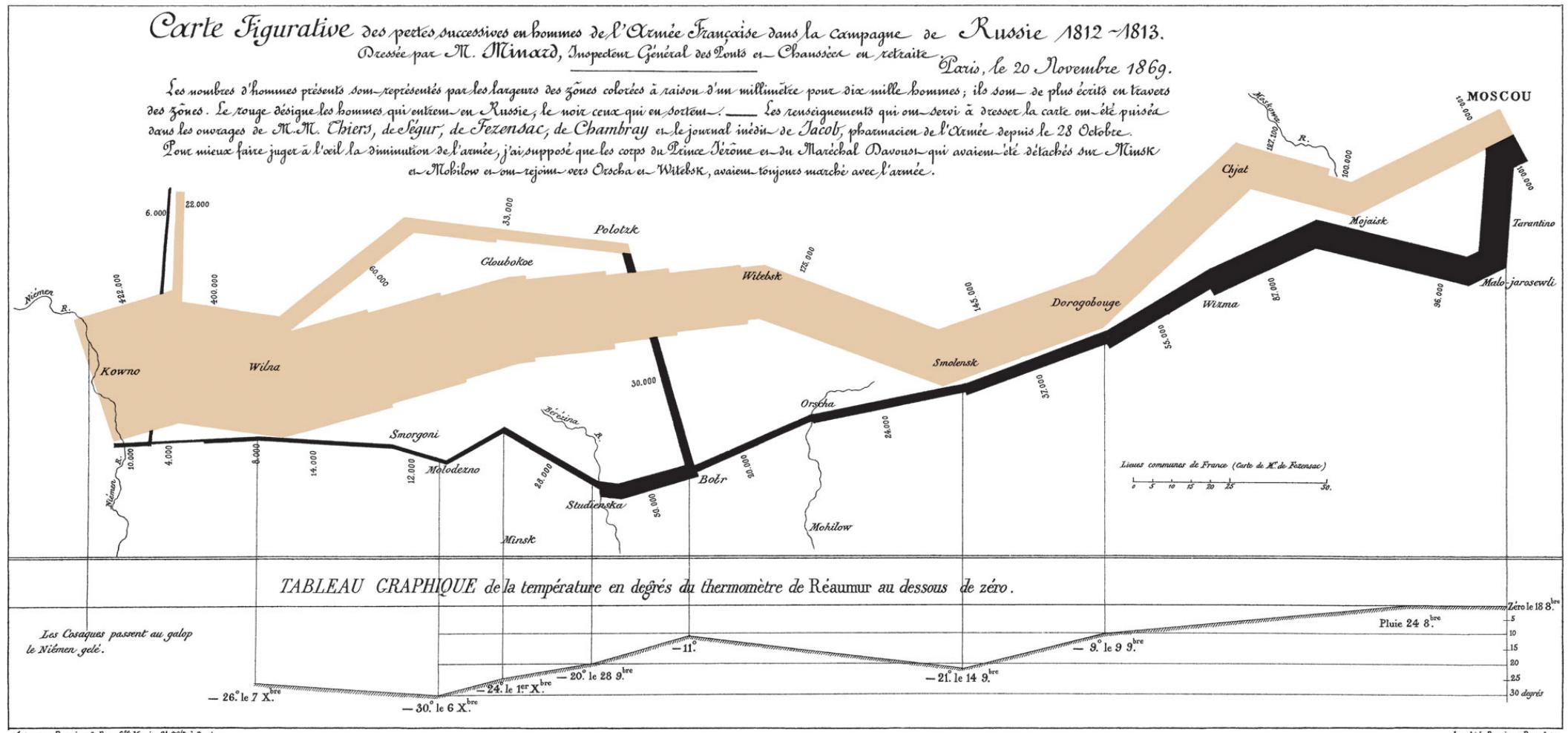
# Advanced Visualization

INTRODUCTION



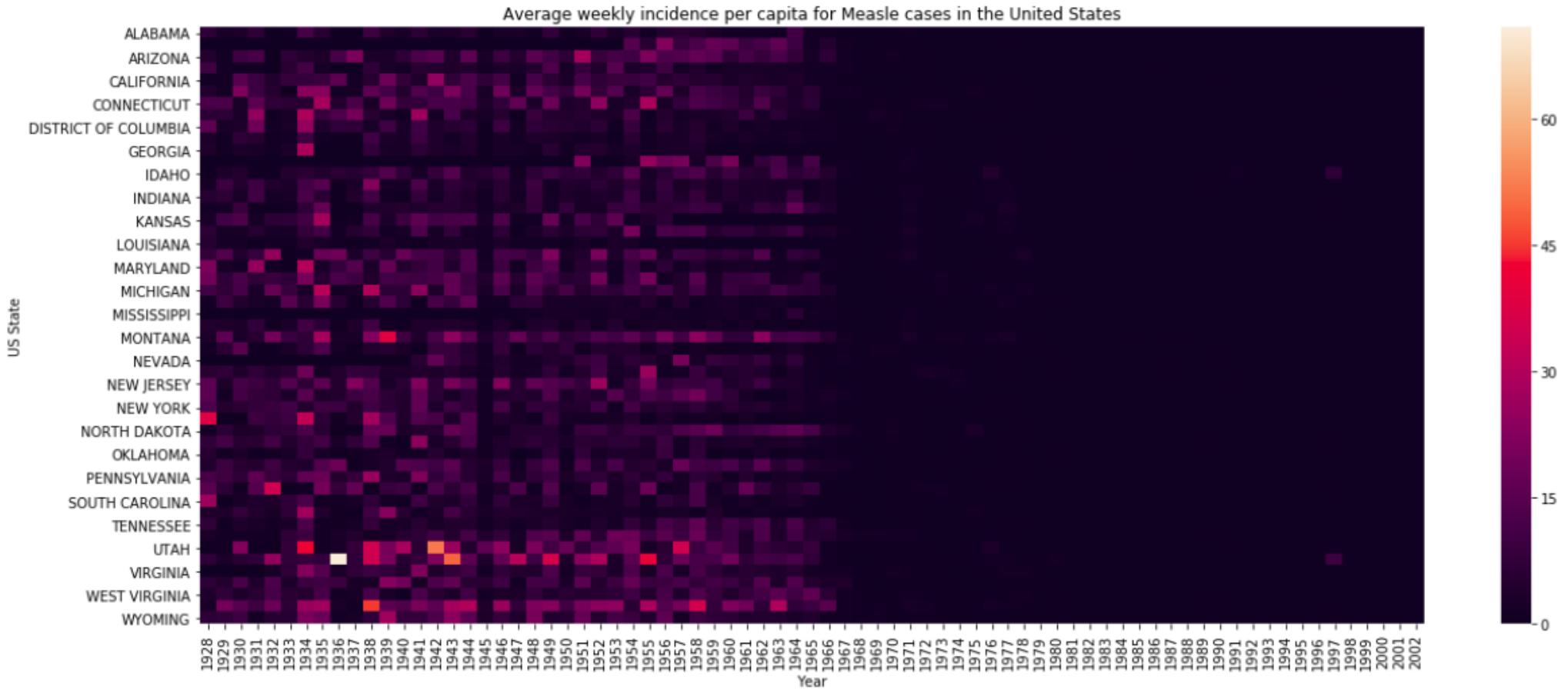
To DATA SCIENCE

# Minard's Napoleon March



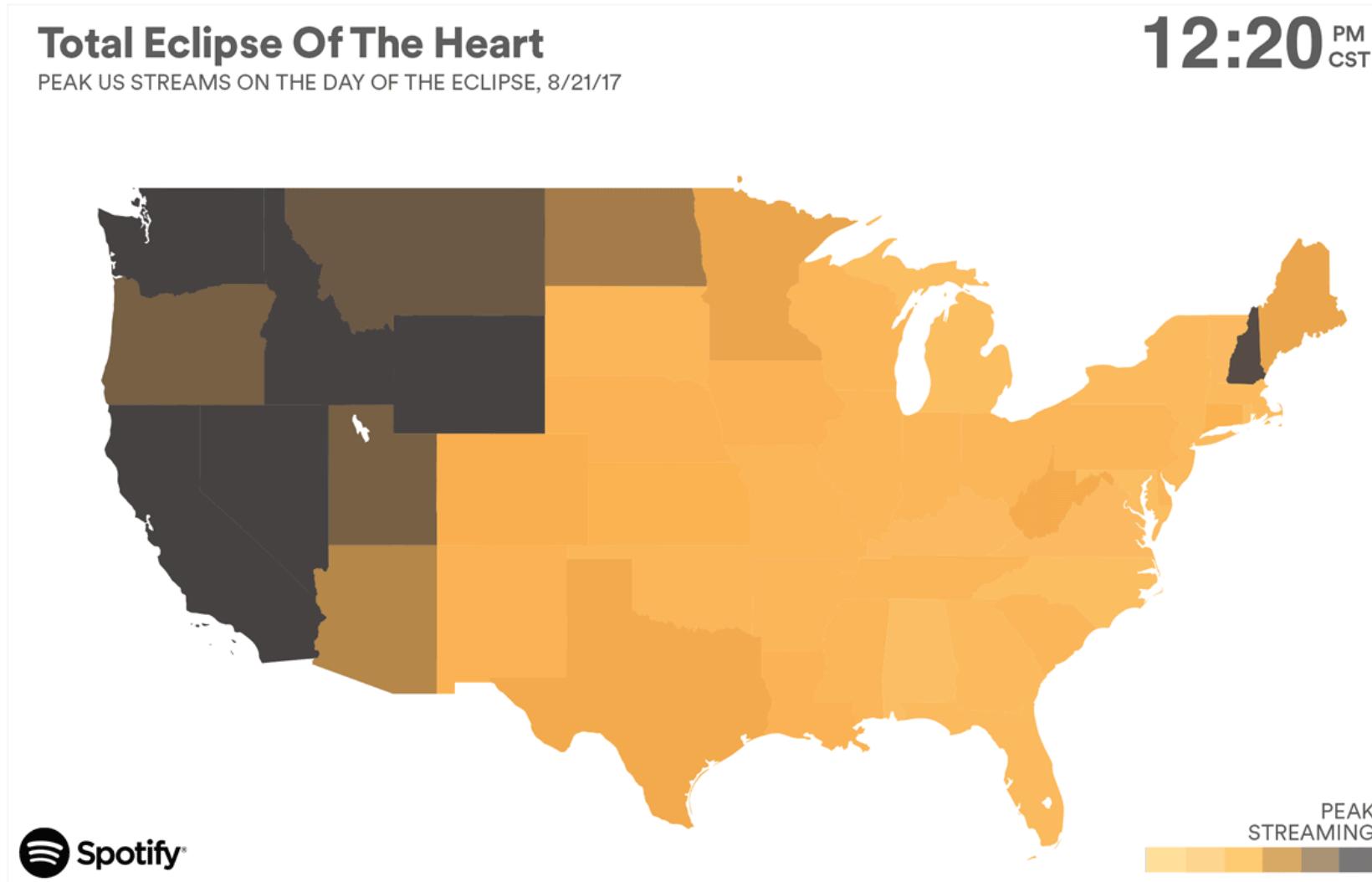
# Heatmaps

## source



# Spotify: Total Eclipse of the Heart

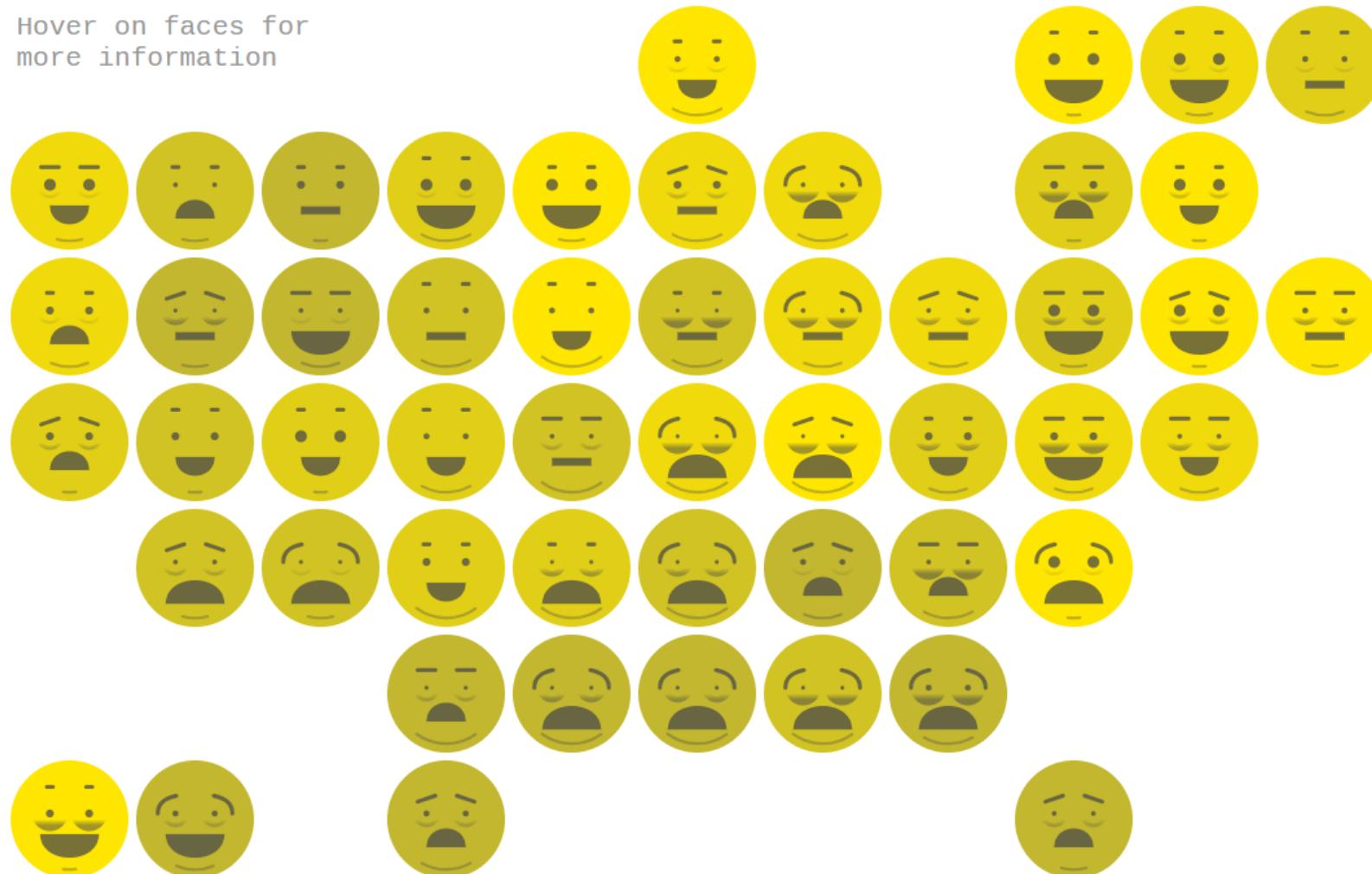
source



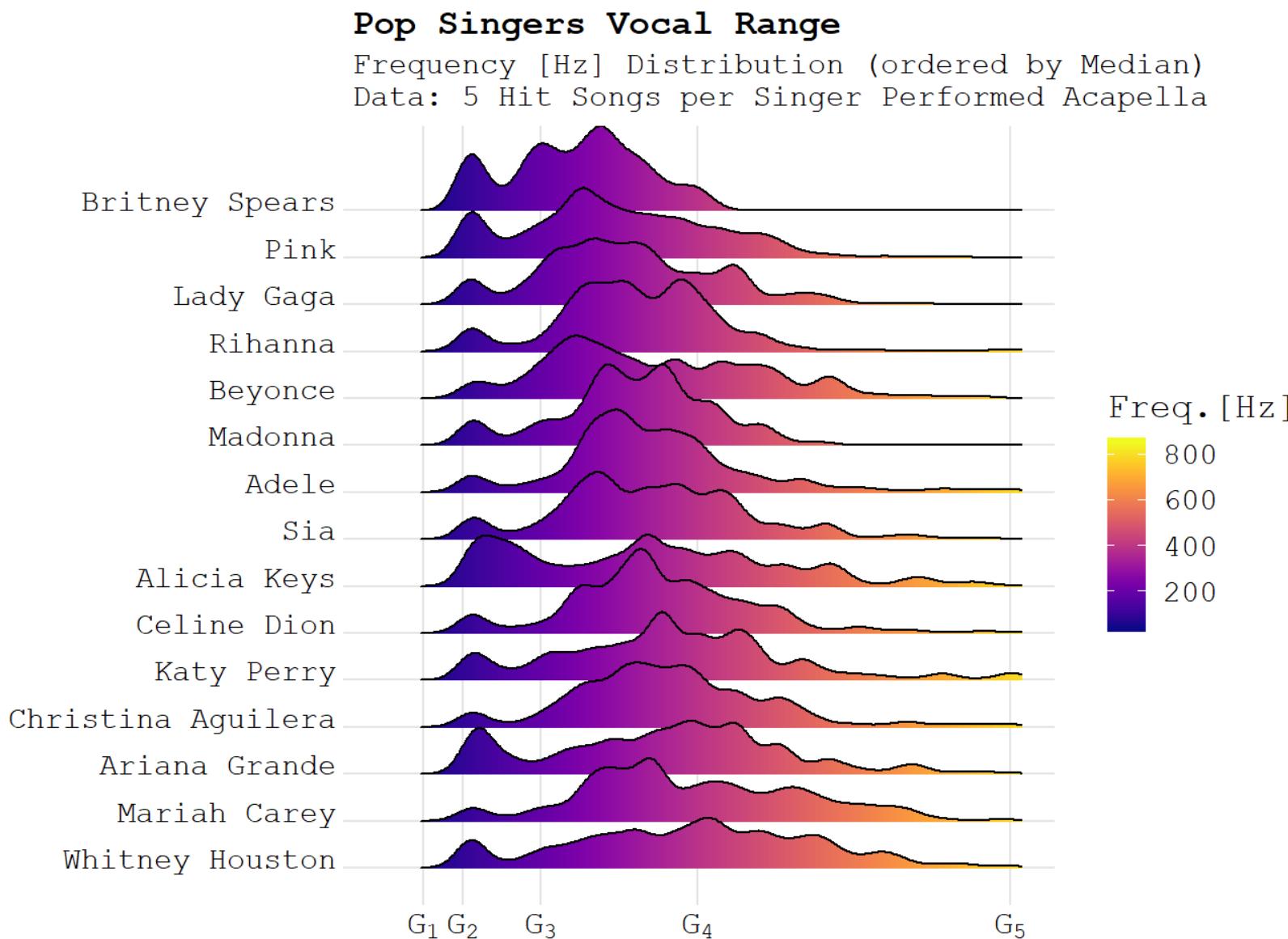
# Chernoff Faces

[source](#)

Hover on faces for  
more information



# Ridge plot (a.k.a Joy plot)



# Flowing Data: A Day in the Life of Americans

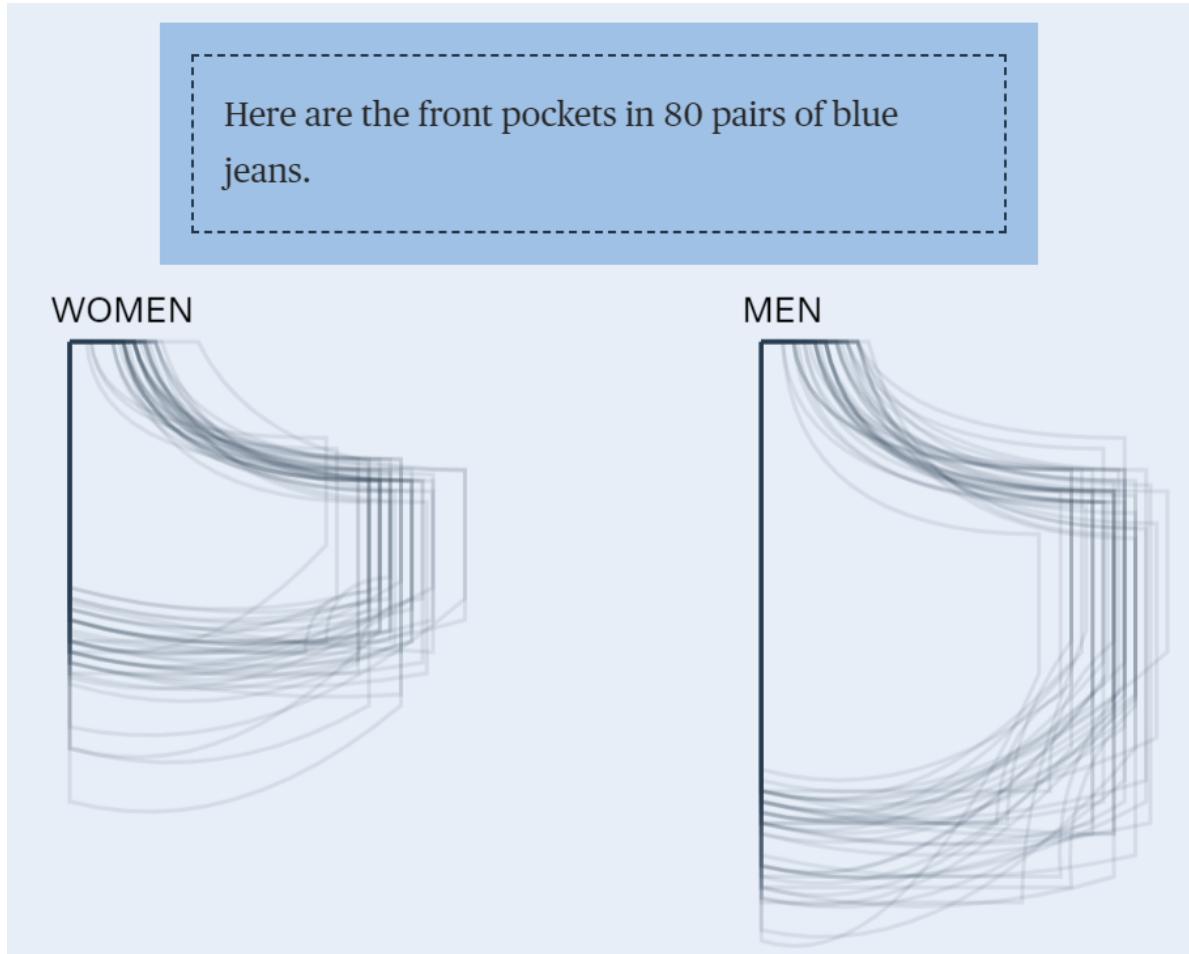
[source](#)

One Day in the Life of Americans by Flowingdata



# Pudding: Women's Pockets (and every single post on their site!)

source



# The Gapminder story

The history of the world encapsulated in a simple visualization:

[source](#)

# Dangers of Dirty Data

INTRODUCTION



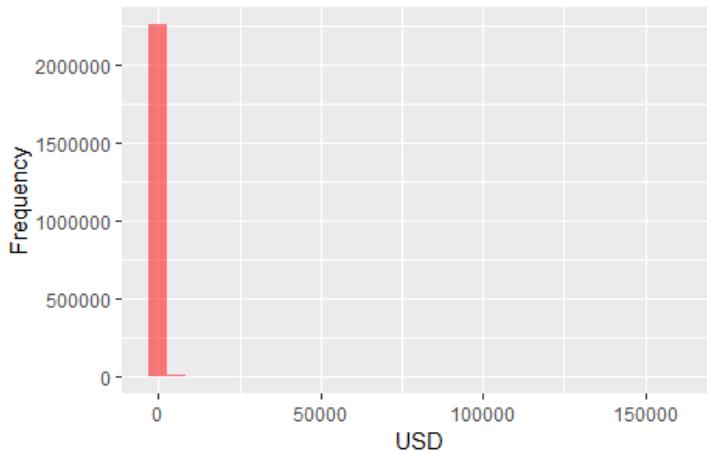
To DATA SCIENCE

# What could be dirty about data?

- The data itself
- The data's structure

# The data itself: Outliers

- Numerical Outliers: This is a histogram of random ~2.3 million transactions on ebay US website in over a few weeks in 2013 ([source](#)):



- Textual Outliers: The [Blog Authorship Corpus](#) consists of 19K posts by bloggers from blogger.com in 2004. These are actual words used in the 10-20 age group:

aaaaaaaaaaaaaaaaaaaaargh, lolzi, jfjgfjhgfjfjf, roflmfao, duuuuuuh, walang, dunno

# The data itself: Missing data

source

כל כתבות היום

## קמיל פוקס על כישלון המדגמים: "כשיפרו לי על התוצאות, חשבתי שאני מת"

בכיר באחד העורכים: "טעינו בגזול. גילינו פער של 5%-6% בתוצאות בין המדגם לבין הקלפי - אנשים נכנסו לקלפי המדגמית ושיקרו" ■ ד"ר מינה צמח: "סגרנו את קלפיות המדגם בשעה שמונה וחצי בערב אבל ראיינו שככל שהשעה יותר מאוחרת, כך הליכוד מתחיזק"



עדי דברת-מורץ | התראות במייל  
14:52 18.03.2015

פרופסור קמיל פוקס שמע לראשונה על כך שתוצאות המדגם שהציג אטמול בערך 10 שניות רק כשהתעורר היום בבוקר משיחת טלפון של מפיק ערוץ 10, שהזמין אותו לאולפן כדי לדון **בתוצאות האמת של הבחירות**. "זונתי אמרה לי 'שמעת מה קרה?' וסיפרה לי על תוצאות האמת, ואני חשבתי שאני מת".

פוקס סיפר כי "זו לא הפעם הראשונה שזה קורה לי - בשנת 96' אמרתי לחיים בין שתוצאות המדגם מלמדות על תיקו עם יתרון קל לפרס ואוז הلتכי לשון עם פרס והתעוררתי עם ביבי (נתניהו). אבל הפעם התחששה הרבה יותר נרואה כי הפער בין תוצאות המדגם לתוצאות האמת יותר גדולות ובשנת 96' הפער היה רק של 16 אלף קולות", אומר פוקס ל-TheMarker ביום לאחריו מה שייזכר



**קריאה Zus**



# The data's structure

In a word: Excel.

source

Australian Bureau of Statistics											
1800.0 Australian Marriage Law Postal Survey, 2017											
Released on 15 November 2017											
Table 5 Participation by Federal Electoral Division(a), Males and Age											
Gender apartheid											
<b>Yeah NA</b>		18-19 years	20-24 years	25-29 years	30-34 years	35-39 years	40-44 years	45-49 years	50-54 years	55-59 years	60-64 years
Lingiari(c)	Total participants	292	1,058	1,465	1,653	1,515	1,516	1,710	1,730	1,753	1,574
Lingiari(c)	Eligible participants	572	2,910	3,789	3,996	3,607	3,506	3,645	3,331	2,960	2,456
Lingiari(c)	Participation rate (%)	51.0	36.4	38.7	41.4	42.0	43.2	46.9	51.9	59.2	64.1
<b>Primary keynotes</b>		<b>Comma on</b>									
Merged cells	Total participants	442	1,461	2,066	2,357	2,188	2,057	2,224	2,108	2,134	1,772
Solomon	Eligible participants	750	2,991	3,994	4,155	3,634	3,398	3,427	3,066	2,931	2,355
Solomon	Participation rate (%)	58.9	48.8	51.7	56.7	60.2	60.5	64.9	68.8	72.8	75.2
Northern Territory (Total)	Total participants	734	2,519	3,531	4,010	3,703	3,573	3,934	3,838	3,887	3,346
Northern Territory (Total)	Eligible participants	1,322	5,901	7,783	8,151	7,241	6,904	7,072	6,397	5,891	4,811
Northern Territory (Total)	Participation rate (%)	55.5	42.7	45.4	49.2	51.1	51.8	55.6	60.0	66.0	69.5
Australian Capital Territory Divisions	<b>Covariate as Subheading</b>										<b>Summary of data inside data</b>
Canberra(d)	Total participants	1,764	4,789	4,817	4,973	4,626	4,453	5,074	4,826	5,169	4,394
Canberra(d)	Eligible participants	2,260	6,471	6,448	6,509	5,983	5,805	6,302	5,902	6,044	5,057
Canberra(d)	Participation rate (%)	78.1	74.0	74.7	76.4	77.3	76.7	80.5	81.8	85.5	86.9
Fenner(e)	Total participants	1,477	4,687	5,178	5,786	6,025	5,463	5,191	4,208	3,948	3,465
Fenner(e)	Eligible participants	1,904	6,354	7,121	7,822	7,960	7,155	6,480	5,206	4,692	3,945
Fenner(e)	Participation rate (%)	77.6	73.8	72.7	74.0	75.7	76.4	80.1	80.8	84.1	87.8
Australian Capital Territory (Total)	<b>NA Yeah</b>	3,241	9,476	9,993	10,755	10,051	9,316	10,265	9,034	9,117	7,059
Australian Capital Territory (Total)	Eligible participants	4,164	12,825	13,569	14,331	13,943	12,960	12,782	11,108	10,736	9,002
Australian Capital Territory (Total)	Participation rate (%)	77.8	73.9	73.7	75.1	76.4	76.5	80.3	81.3	84.9	87.3
Australia	Total participants	151,297	438,166	441,658	460,548	462,206	479,360	524,620	517,693	543,449	506,799
Total	Eligible participants	201,439	635,909	646,916	665,250	656,446	660,841	693,850	659,150	664,720	597,386
Total	Participation rate (%)	75.1	68.9	68.3	69.2	70.4	72.5	75.6	78.5	81.8	84.8
(a) The Federal Electoral Divisions are current as at 24 August 2017											
(b) Includes those whose age is unknown											
(c) Includes Christmas Island and the Cocos (Keeling) Islands											
(d) Includes Norfolk Island											
(e) Includes Jervis Bay											
<b>Return of the table junk</b>											
<b>MS Excel or Die</b>											

# Spreadsheet Blunder

source

BBC [Sign in](#) | Home News Sport Reel Worklife Travel

## NEWS

Home | War in Ukraine | Climate | Video | World | UK | Business | [Tech](#) | Science | Entertainment & Arts

Tech

### Excel: Why using Microsoft's tool caused Covid-19 results to be lost

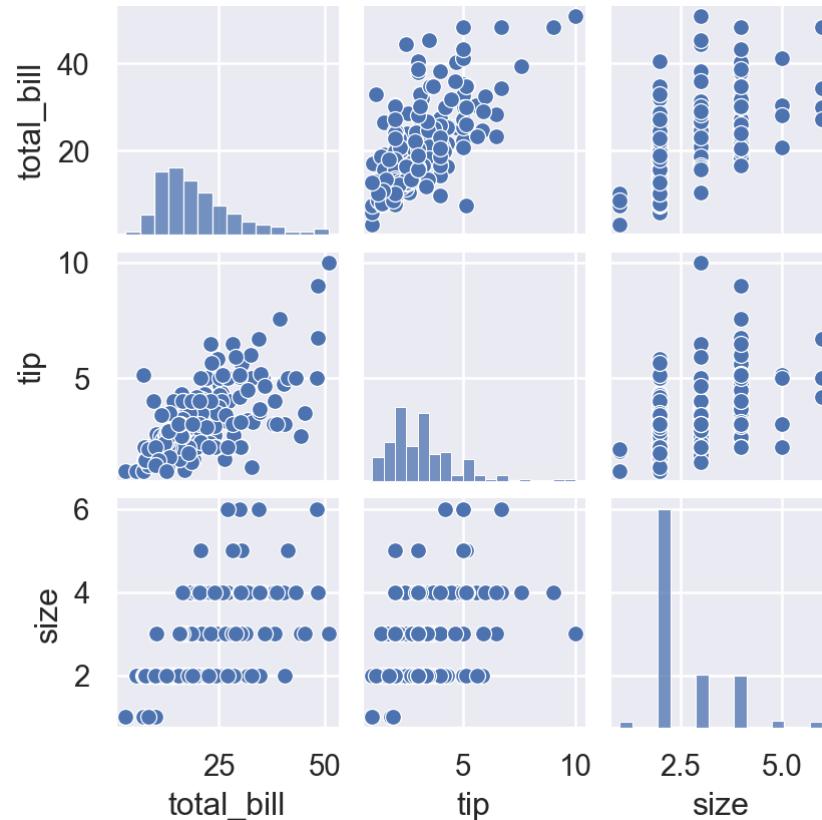
⌚ 5 October 2020



# Some advice on cleaning data

Plot first:

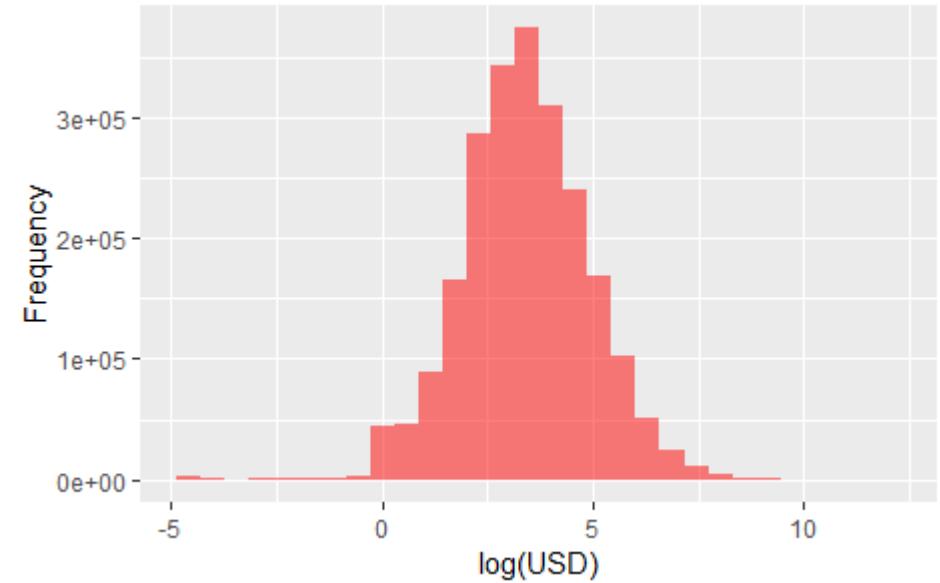
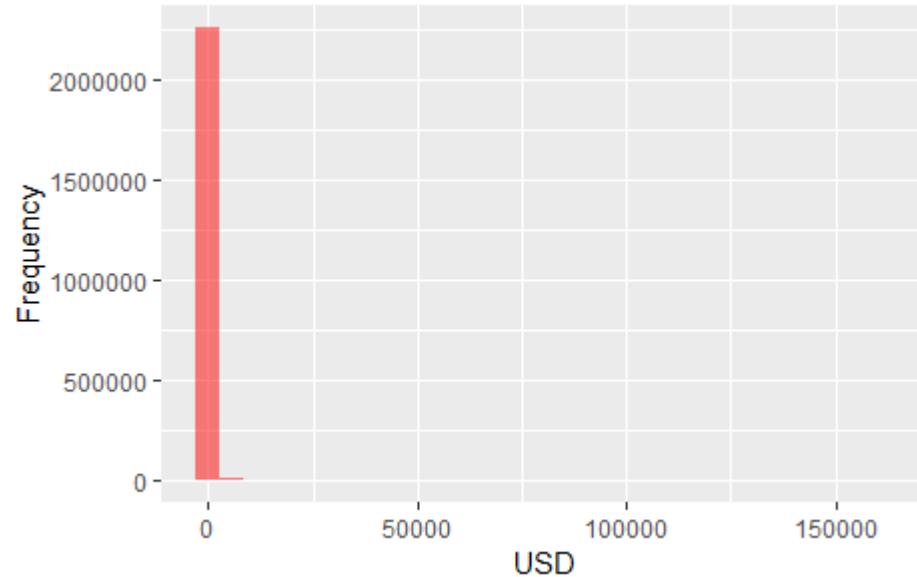
```
tips = sns.load_dataset('tips')
sns.pairplot(tips, height=1.5)
plt.show()
```



# Some advice on cleaning data

Apply common transformations:

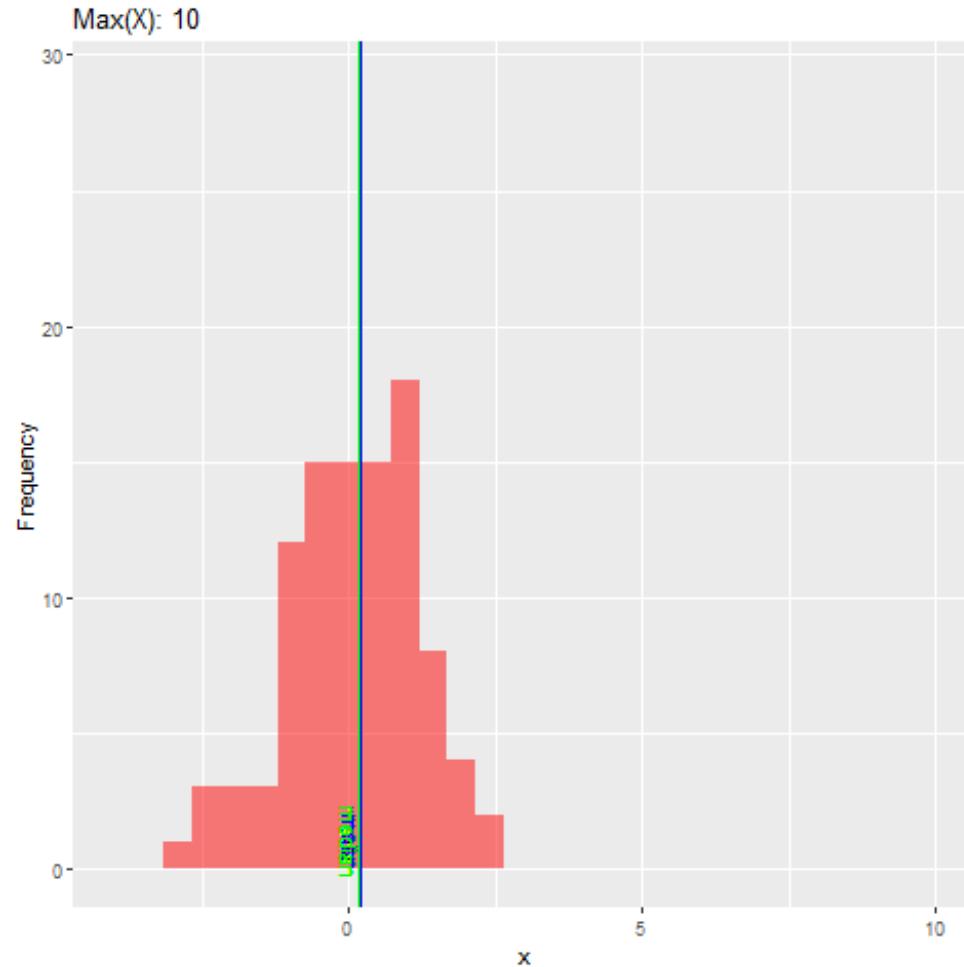
Here's how ebay's 2.3 million transactions look with a log transformation:



# Some advice on cleaning data

Use robust statistics (an entire field in Statistics):

For example the Median is much more robust to extreme values than the mean:



# Some advice on cleaning data

Tidy your data:

Each variable is a column, each observation is a row, and each type of observational unit is a table. ([Hadley Wickham](#))

Untidy (wide):

	SubjectName	Weight_Age20	Weight_Age30
1	John Doe	80	90
2	Jane Doe	60	50

Tidy (long):

	SubjectName	Age	Weight
1	John Doe	20	80
2	John Doe	30	90
3	Jane Doe	20	60
4	Jane Doe	30	50

# Intro to Data Science

INTRODUCTION



To DATA SCIENCE