# **Doing Visual Data Science**

# Foundations, Techniques and Practice

Cagatay Turkay

Professor,

Centre for Interdisciplinary Methodologies

University of Warwick

CIM WARWICK

# Next up .. VDS DIY

**VDS DIY BLOCK #1**

**Thinking Visually**

**Part-1: Re/De-constructing Visualisations**
- Introduction to visualisation basics
- Moving on to the visualisation grammar
- A "language" of visualisation
- Altair -- a versatile visualisation library in Python

**Part-2: Foundational VDS Techniques**
- Visualisation for interrogating the data
  - Conformity, Outliers, Shapes
- Multiple perspectives
  - Using small multiples

**VDS DIY BLOCK #2**

**Visual Data Science -- Getting "involved"**

**Part-1: Working with models**
- Working with models visually
  - Interpreting clustering
  - Working with projections for high-dimensional data
- Beyond accuracy

**Part-2: Making it interactive**
- Bringing in interaction
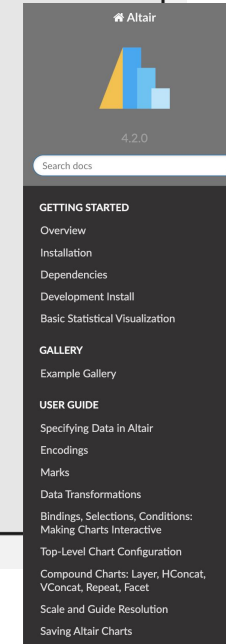- Linking it up
- Going back to the data

**VDS DIY BLOCK #3**

**Doing (V)DS responsibly**

**Part-1: Being aware**
- Dance of the p-values
- Maps and colouring

**Part-2: Communicating openly**
- Thinking about the narrative
- Communicating responsibly



- Hands-on exercises to explore some of these practices/ideas
- Will use Altair in Jupyter Notebooks to visualise and bring interactivity

# Doing **Visual Data Science**

# Foundations, Techniques and Practice

**DIY BLOCK #1 - Thinking Visually**

Part-1: Re/De-constructing Visualisations

# BUILDING BLOCKS OF VISUALISATIONS
*VISUAL ELEMENTS*

# Marks & Channels

**Visual Marks**: represent **items (**or **links)**

**Channels**: change **appearance** based on **attribute**

**Channel = Visual Variable**

# Basic visual elements

**VISUAL MARKS**
i.e., primitives

→ Points    → Lines    → Areas

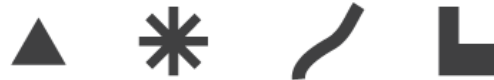---

→ Position

→ Horizontal    → Vertical    → Both

→ Color

**VISUAL CHANNELS**
i.e., how they look

→ Shape

→ Tilt

→ Size

→ Length    → Area    → Volume

From Visualization Analysis & Design by T. Munzner, 2014

# Visual encoding

**Def.** *Representing information visually as combination of marks and channels*

1:
vertical position

mark: line

2:
vertical position
horizontal position

mark: point

3:
vertical position
horizontal position
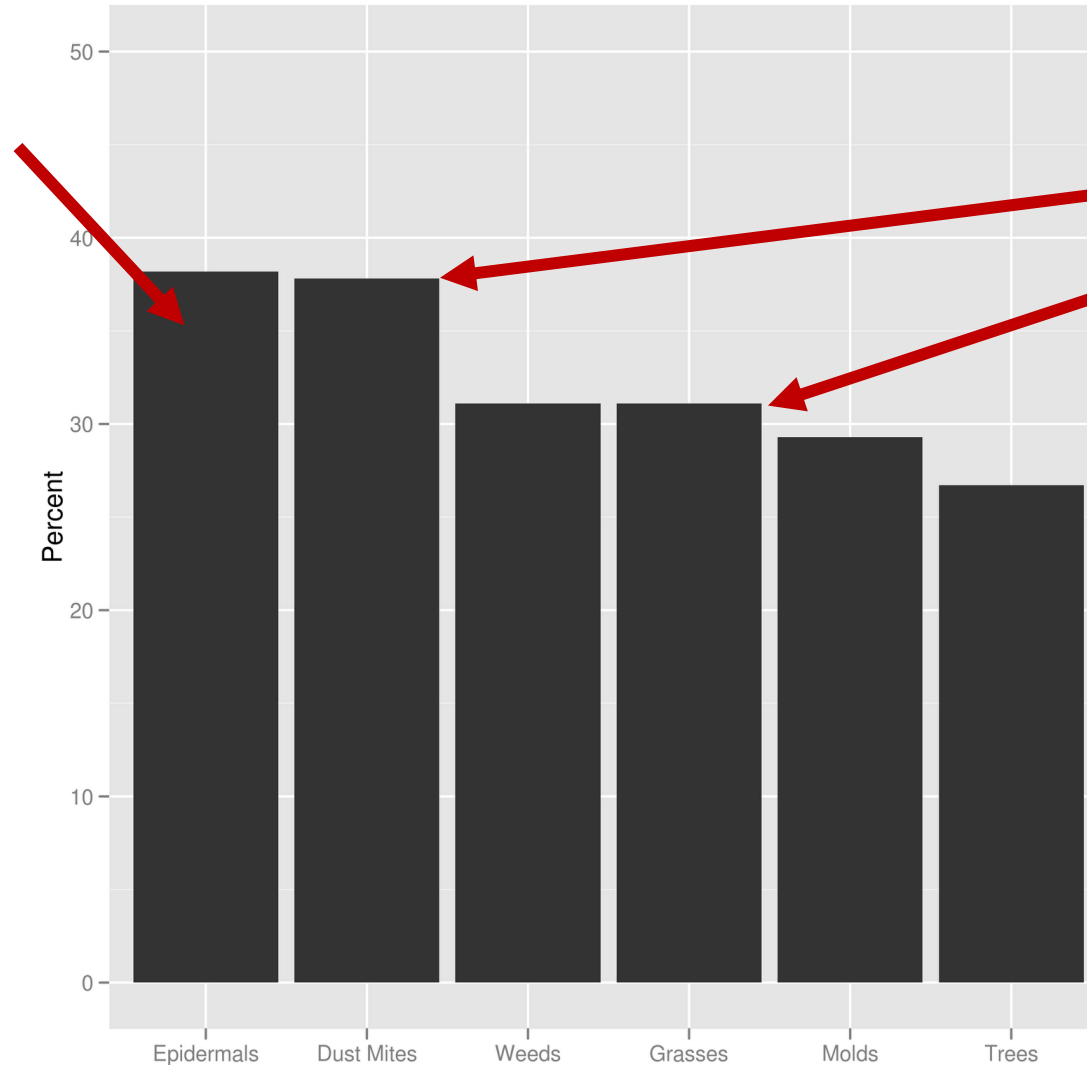color hue

mark: point

4:
vertical position
horizontal position
color hue
size (area)

mark: point

# Basic visual elements in action



VISUAL MARK
lines (**thick** ones)

VISUAL CHANNEL
length (of bars)

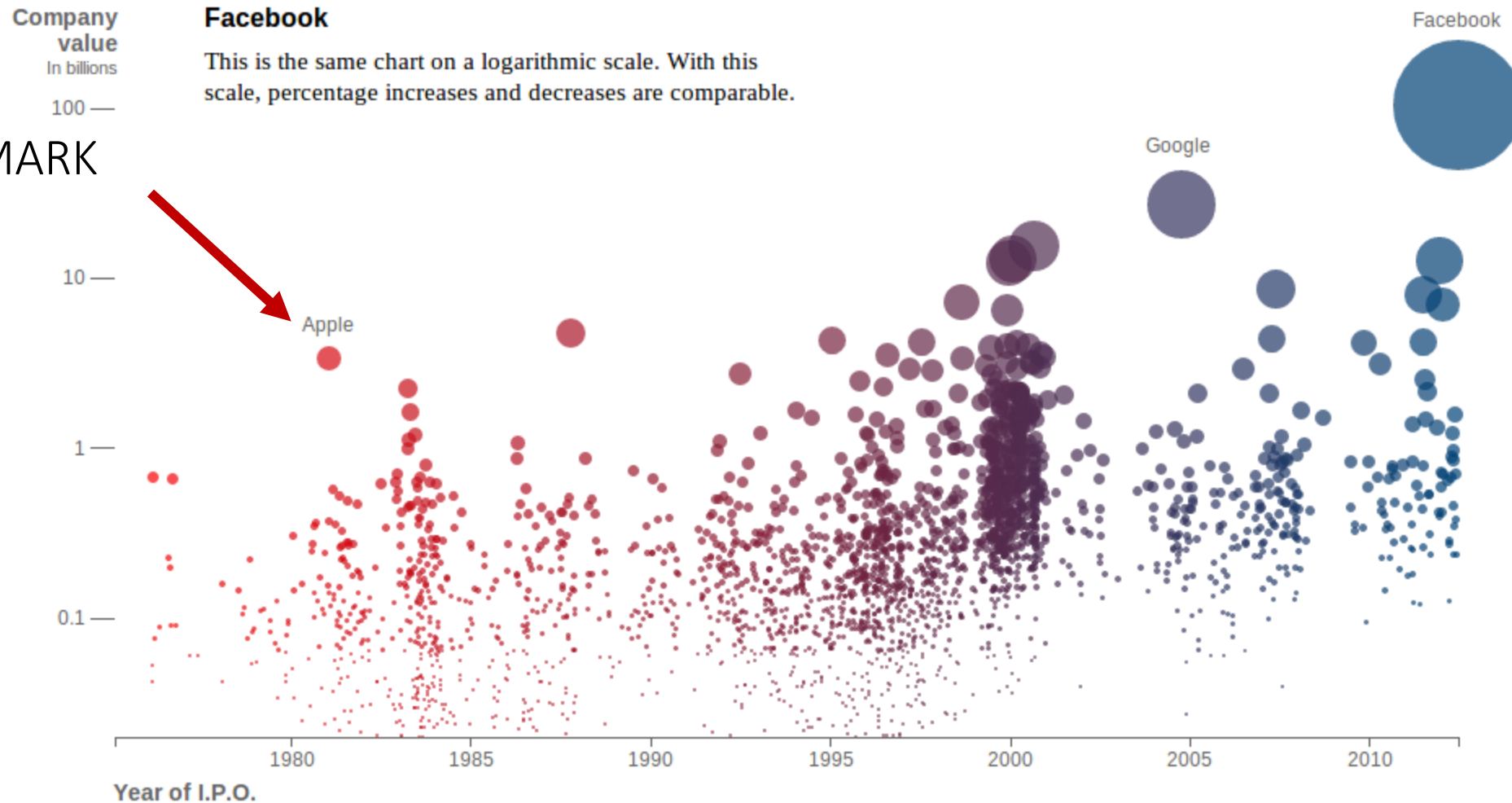http://rrubyperlundich.blogspot.co.uk/2012_09_01_archive.html

# Basic visual elements in action

VISUAL CHANNELS
position x and y
colour
Area (redundant)

VISUAL MARK
points

**Company value**
In billions

**Facebook**

This is the same chart on a logarithmic scale. With this scale, percentage increases and decreases are comparable.



Facebook

Google

Apple

100 —

10 —

1 —

0.1 —

1980        1985        1990        1995        2000        2005        2010

**Year of I.P.O.**

http://www.nytimes.com/interactive/2012/05/17/business/dealbook/how-the-facebook-offering-compares.html

# FIND THE ODD ONE OUT!

Length

Size

Orientation

Lightness

Hue

# Visual channels don't work equally well!

# A Layered Grammar of Graphics

Hadley WICKHAM

| GPL | ggplot2 |
|---|---|
| DATA ⟶ | Defaults |
| TRANS | Data Mapping |
| ELEMENT ⟶ | Layer |
| | Data Mapping Geom Stat Position |
| SCALE ⟶ | Scale |
| GUIDE ⟶ | |
| COORD ⟶ | Coord |
| | Facet |

Statistical graphic specifications are expressed in six statements:

1) DATA: a set of data operations that create variables from datasets,
2) TRANS: variable transformations (*e.g.*, *rank*),
3) SCALE: scale transformations (*e.g.*, *log*),
4) COORD: a coordinate system (*e.g.*, *polar*),
5) ELEMENT: graphs (*e.g.*, *points*) and their aesthetic attributes (*e.g.*, *color*),
6) GUIDE: one or more guides (*axes*, *legends*, etc.).

# **Vega-Lite** – A Grammar of Interactive Graphics



**Vega-Lite** is a high-level grammar of interactive graphics. It provides a concise, declarative JSON syntax to create an expressive range of visualizations for data analysis and presentation.

Vega-Lite specifications describe visualizations as encoding mappings from data to **properties of graphical marks** (e.g., points or bars). The Vega-Lite compiler **automatically produces visualization components** including axes, legends, and scales. It determines default properties of these components based on a set of **carefully designed rules**. This approach allows Vega-Lite specifications to be concise for quick visualization authoring, while giving user control to override defaults and customize vari⟨…⟩ As we also designed Vega-Lite to support data analysis, Vega-Li⟨…⟩ **transformations** (e.g., aggregation, binning, filtering, sorting) an⟨…⟩ (e.g., stacking and faceting). Moreover, Vega-Lite specifications ⟨…⟩ layered and multi-view displays, and made **interactive with sele⟨…⟩**

Compared to Vega, Vega-Lite provides a more concise and conv⟨…⟩ specifications to Vega specifications, users may use Vega-Lite a⟨…⟩ level Vega for advanced use cases.

For more information, read our introduction article to Vega-Lite v⟨…⟩ Lite v2, see the documentation and take a look at our example g⟨…⟩

**Get started**
Latest Version: 5.2.0

# Vega-Lite: A Grammar of Interactive Graphics

## Arvind Satyanarayan, Dominik Moritz, Kanit Wongsuphasawat, and Jeffrey Heer
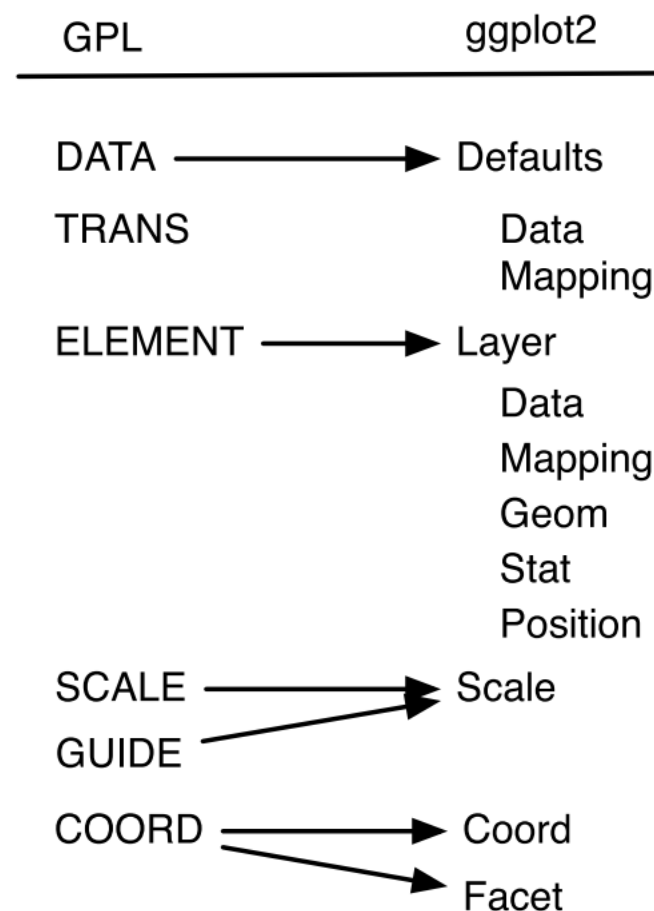


Fig. 1. Example visualizations authored with Vega-Lite. From left-to-right: layered line chart combining raw and average values, dual-axis layered bar and line chart, brushing and linking in a scatterplot matrix, layered cross-filtering, and an interactive index chart.
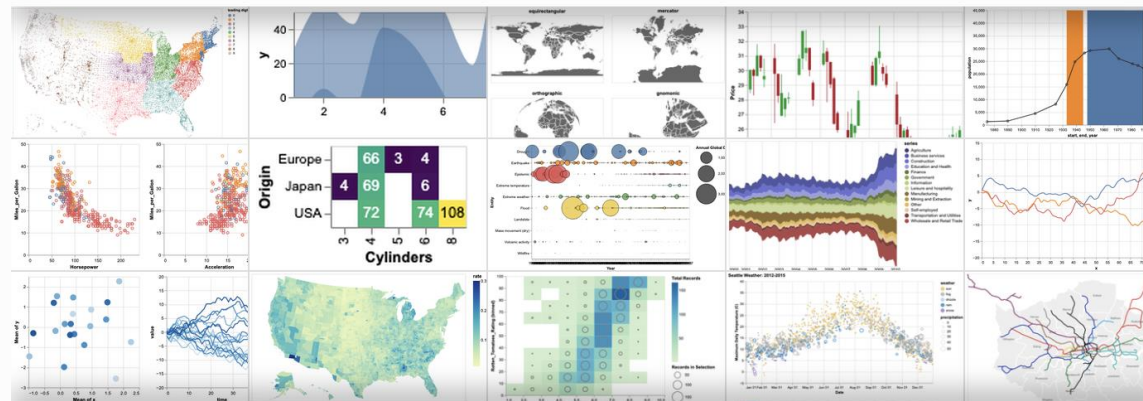
# Altair: Declarative Visualization in Python



Altair is a declarative statistical visualization library for Python, based on Vega and Vega-Lite, and the source is available on GitHub.

With Altair, you can spend more time understanding your data and its meaning. Altair's API is simple, friendly and consistent and built on top of the powerful Vega-Lite visualization grammar. This elegant simplicity produces beautiful and effective visualizations with a minimal amount of code.

## Getting Started

- Overview
- Installation
- Dependencies
- Development Install
- Basic Statistical Visualization

## Gallery

- Example Gallery

# 🏠 Altair

Search docs

# Altair: Declarative Visualization in Python



Altair is a declarative statistical visualization library for Python, based on Vega and Vega-Lite, and the source is available on GitHub.

With Altair, you can spend more time understanding your data and its meaning. Altair's API is simple, friendly and consistent and built on top of the powerful Vega-Lite visualization grammar. This elegant simplicity produces beautiful and effective visualizations with a minimal amount of code.
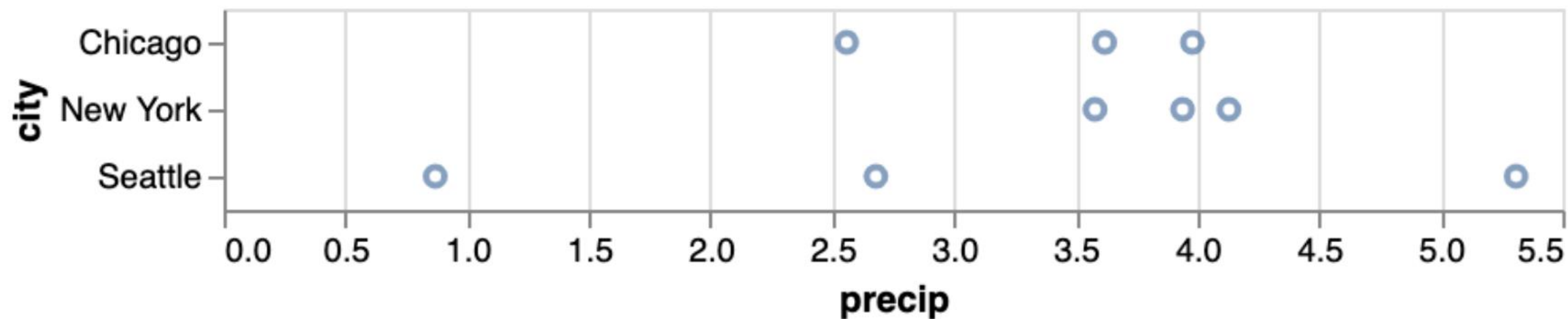
## Getting Started

- Overview
- Installation
- Dependencies
- Development Install
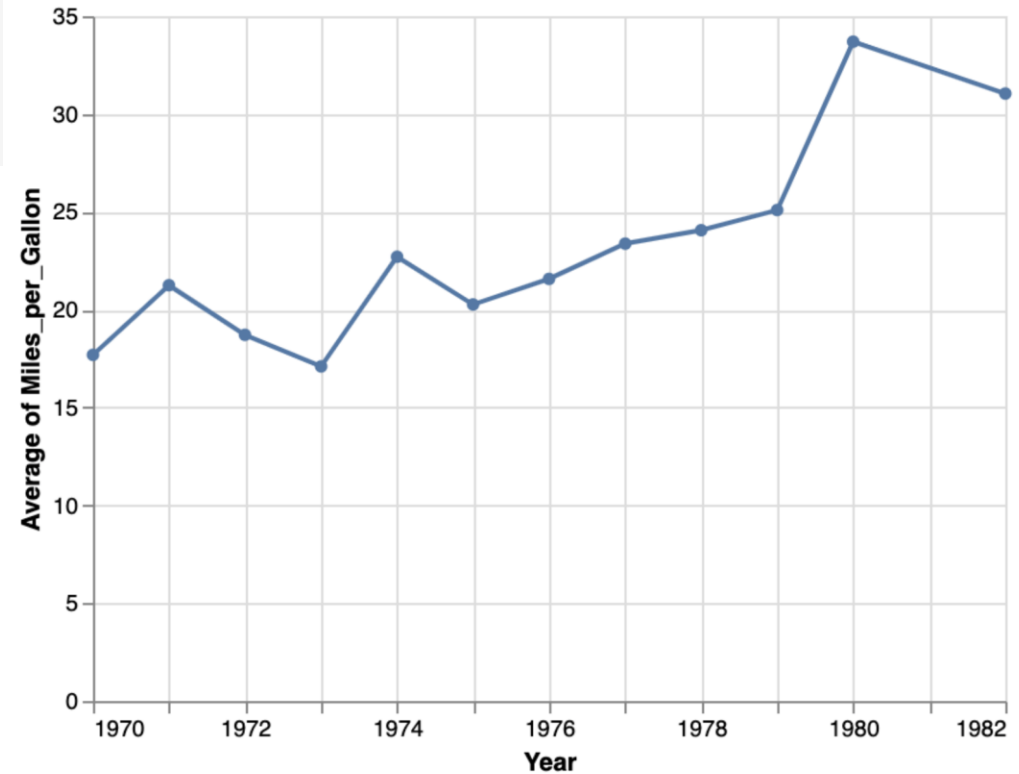- Basic Statistical Visualization

## Gallery

- Example Gallery

```python
alt.Chart(df).mark_point().encode(
    alt.X('precip'),
    alt.Y('city')
)
```

```python
line = alt.Chart(cars).mark_line().encode(
    alt.X('Year'),
    alt.Y('average(Miles_per_Gallon)')
)

point = alt.Chart(cars).mark_circle().encode(
    alt.X('Year'),
    alt.Y('average(Miles_per_Gallon)')
)

line + point
```
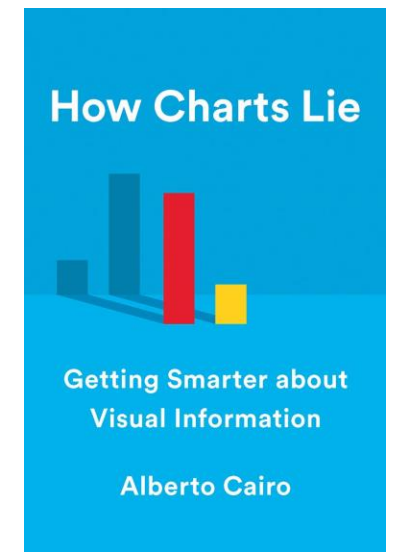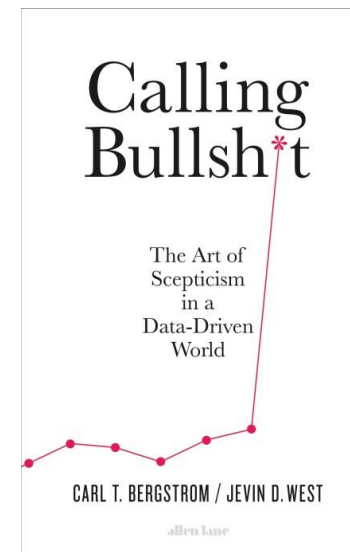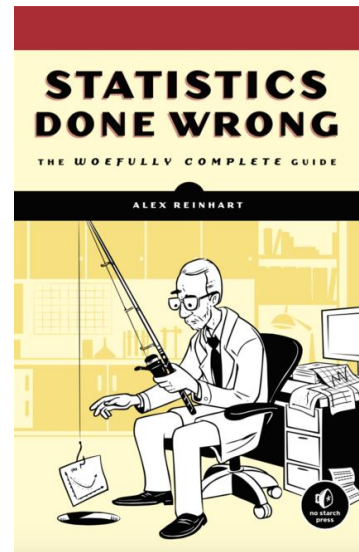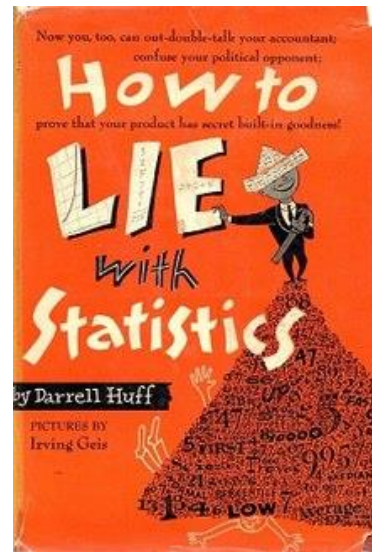
# Doing **Visual Data Science**

# Foundations, Techniques and Practice
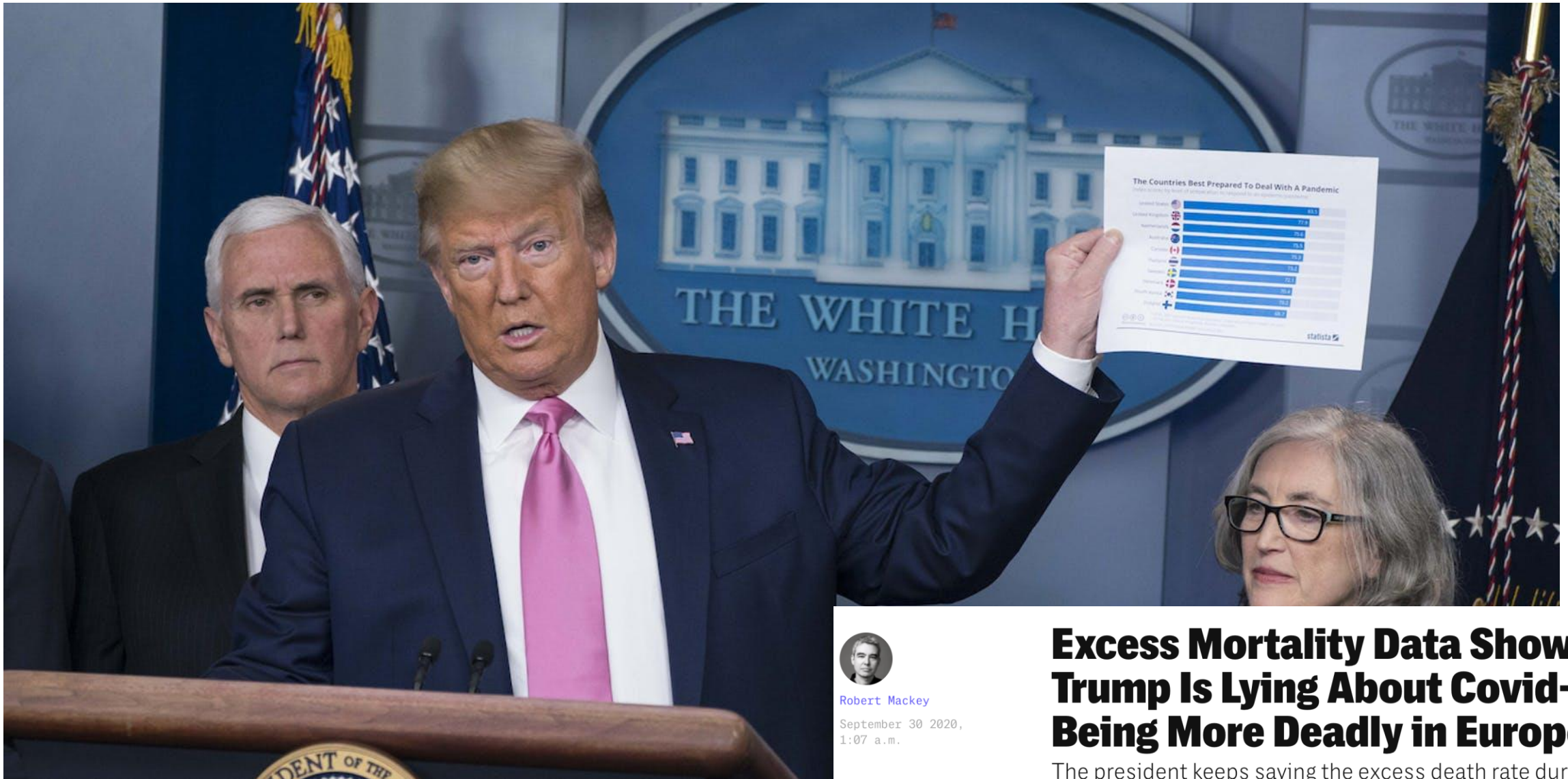
**DIY BLOCK #3 - Doing (V)DS responsibly**

# "There are three kinds of lies: lies, damned lies, and statistics."

(late 1800s, source debated but popularised by Mark Twain, see: https://en.wikipedia.org/wiki/Lies,_damned_lies,_and_statistics )

**Excess Mortality Data Shows Trump Is Lying About Covid-19 Being More Deadly in Europe**

Robert Mackey
September 30 2020,
1:07 a.m.

The president keeps saying the excess death rate during the pandemic is higher in Europe than in the U.S. That's not true.
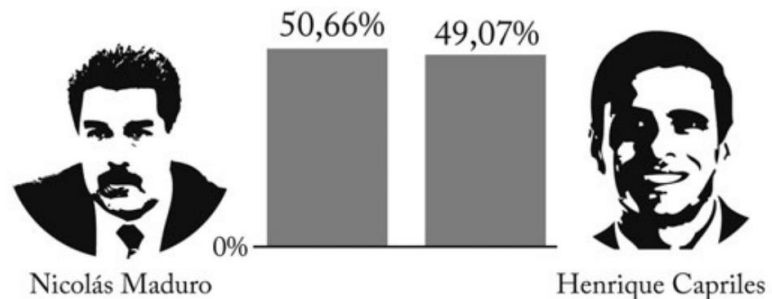
https://theintercept.com/2020/09/29/excess-mortality-data-shows-trump-lying-covid-deadly-europe/

# Truncated axes



PRESIDENTIAL ELECTIONS, 2013

Nicolás Maduro — 50,66%
49,07% — Henrique Capriles

Fig. 5.3 Presidential election results in Venezuela, based on a graphic by Venezonala de Televisión. Notice the truncated Y-axis which greatly distorts the difference between the percentages of vote
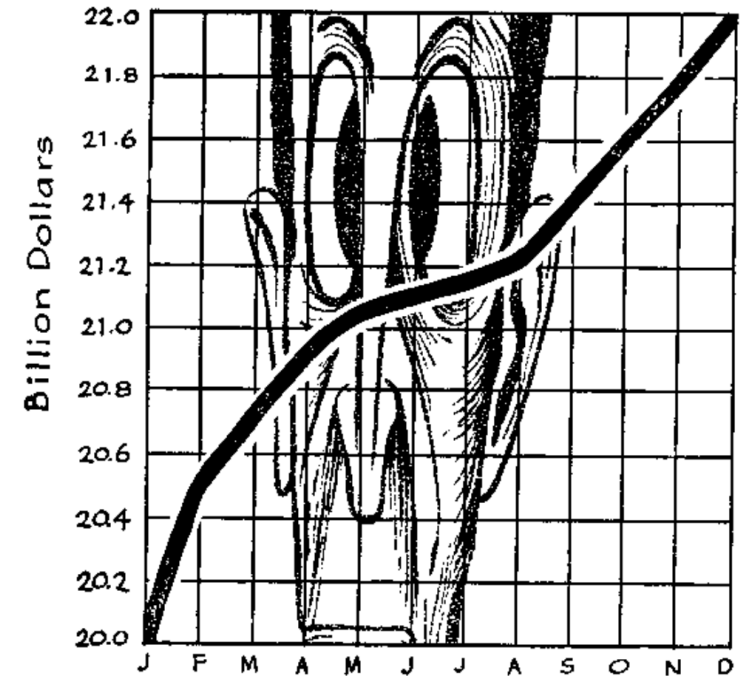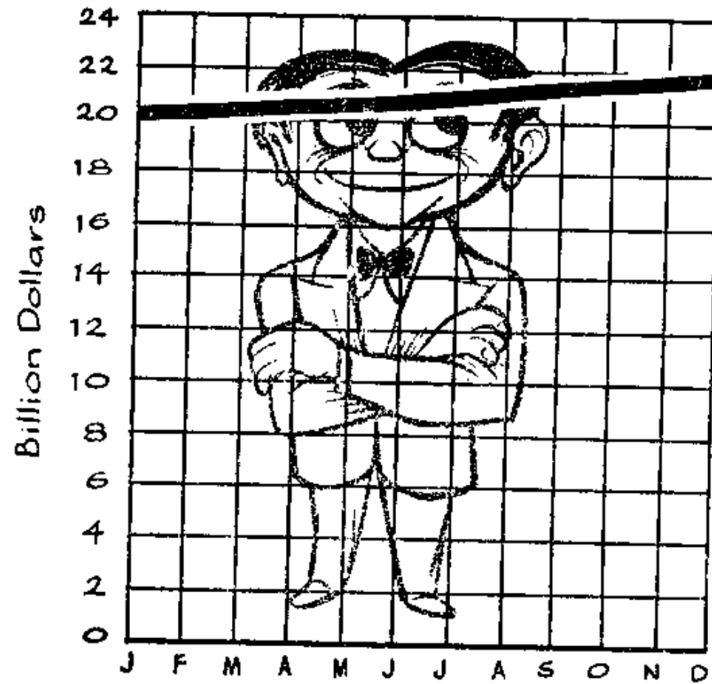
PRESIDENTIAL ELECTIONS, 2013

50,66%   49,07%

Nicolás Maduro     0%     Henrique Capriles

Fig. 5.4 An alternative version of the previous graphic in which a 0-baseline has been added, and the 3D effect has been removed



WELFARE VS. FULL TIME JOBS

108.6M

101.7M

PEOPLE ON WELFARE     PEOPLE WITH A FULL TIME JOB

SOURCE: CENSUS BUREAU, 2011

FOX NEWS FOX NEWS channel

RED SOX BEAT ST. LOUIS CARDINALS 4-2 TO EVEN WORLD SERIES AT T

mediamatters.org

**Chapter 5
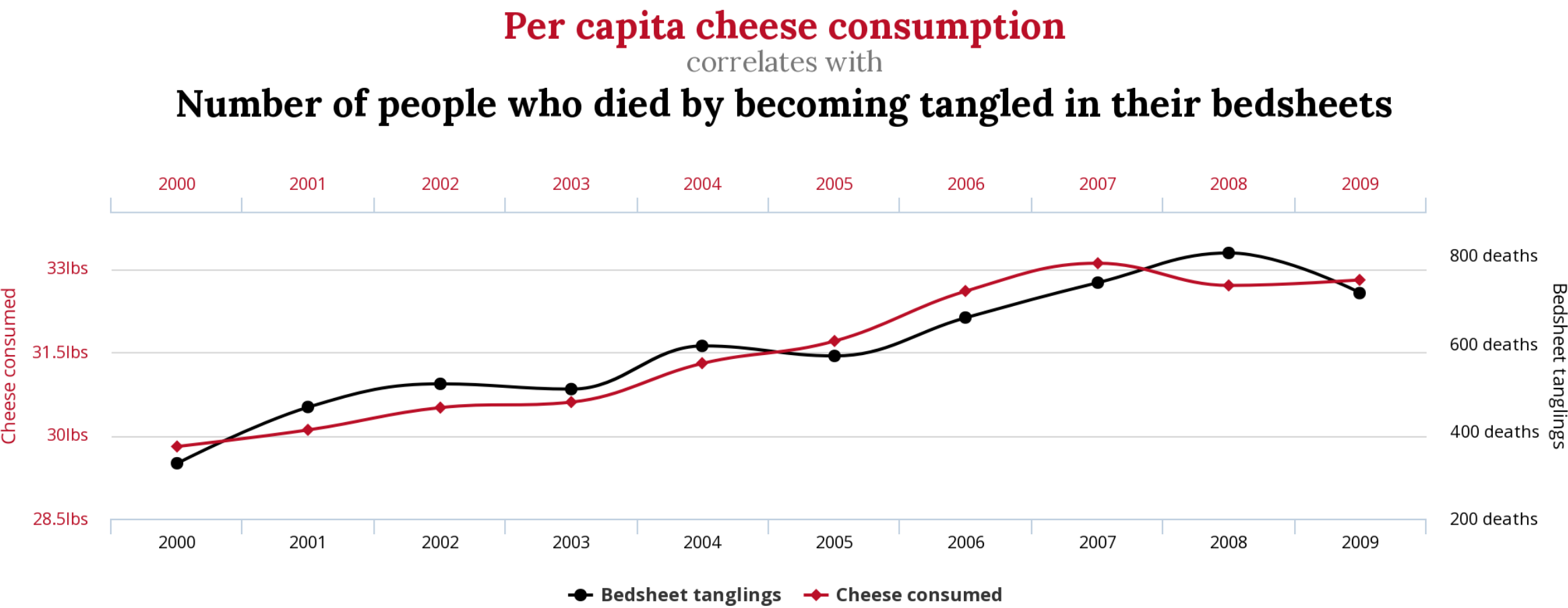Graphics Lies, Misleading Visuals**

**Reflections on the Challenges and Pitfalls
of Evidence-Driven Visual Communication**

Alberto Cairo

# Framing and the power of axes



Figures from: Huff, D., 1993. *How to lie with statistics*. WW Norton & Company.

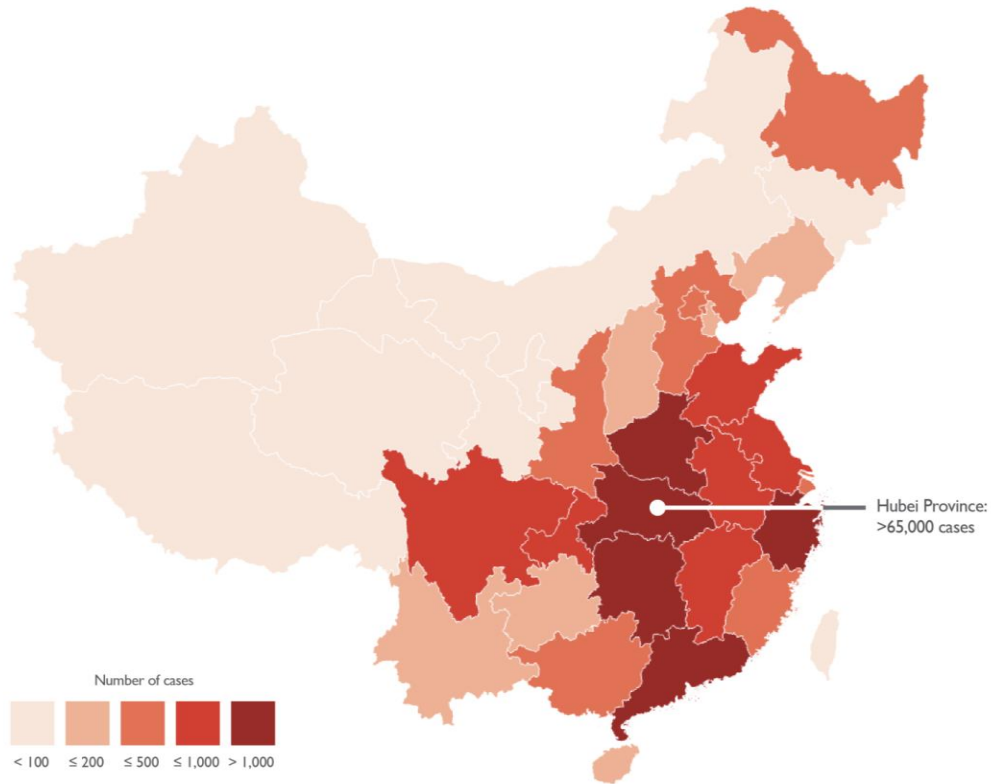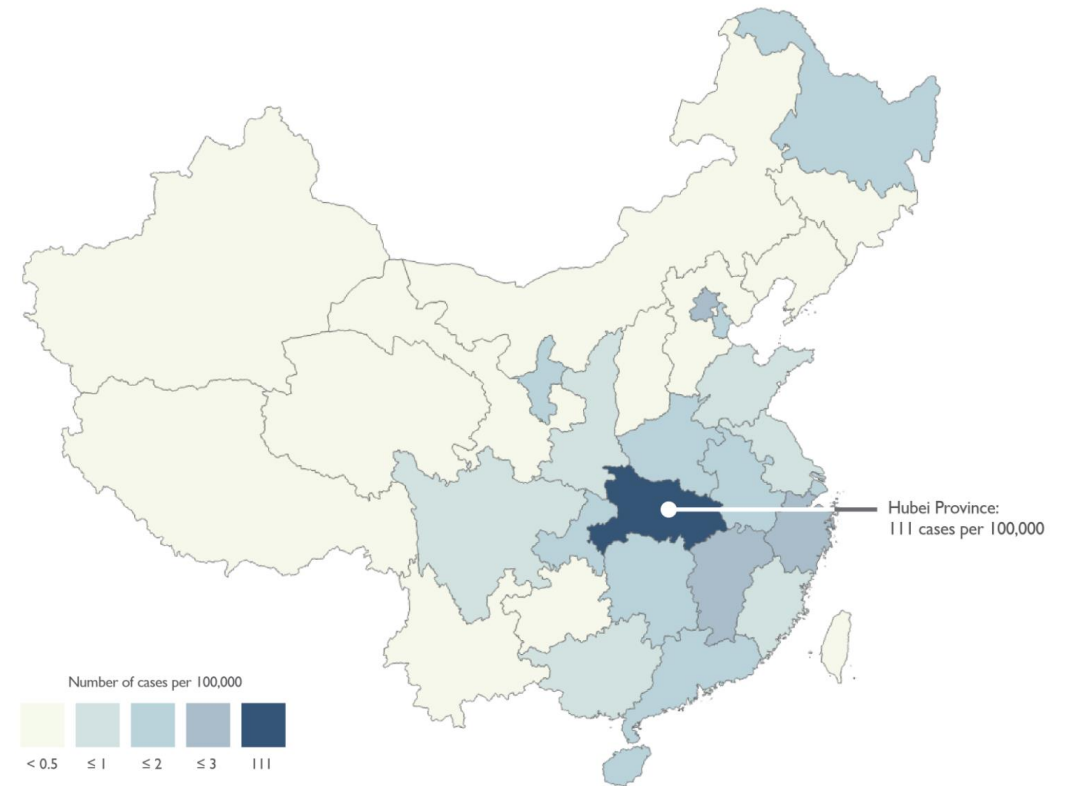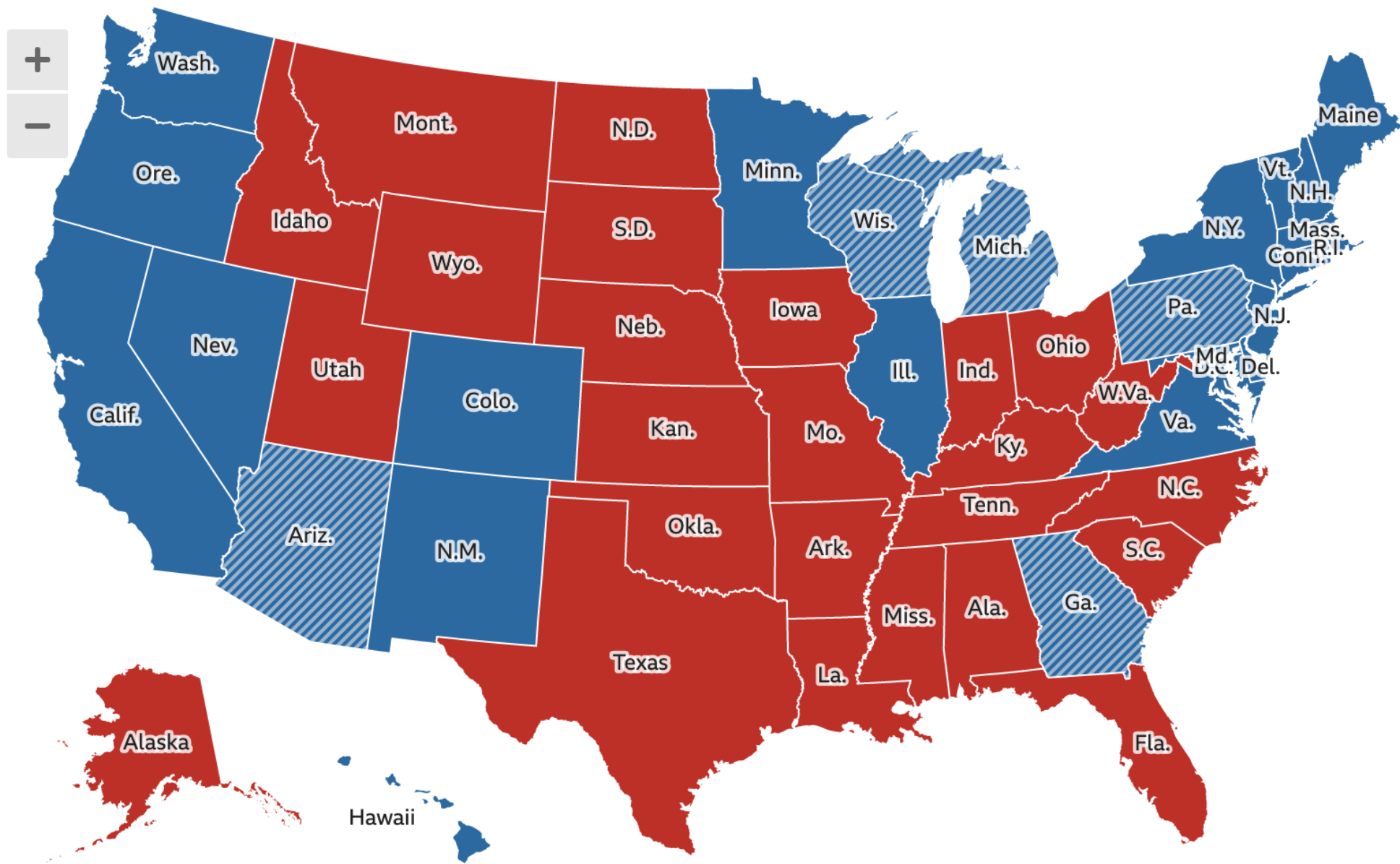# Two axes charts and a false sense of association



**Per capita cheese consumption**
correlates with
**Number of people who died by becoming tangled in their bedsheets**

Bedsheet tanglings    Cheese consumed

tylervigen.com

Spurious correlations

# Absolute/ratios – unfair comparisons, data semantics



Coronavirus in China: 24th February 2020

Hubei Province: >65,000 cases

Number of cases
< 100   ≤ 200   ≤ 500   ≤ 1,000   > 1,000

Coronavirus in China: 24th February 2020

Hubei Province: 111 cases per 100,000

Number of cases per 100,000
< 0.5   ≤ 1   ≤ 2   ≤ 3   111

https://www.esri.com/arcgis-blog/products/product/mapping/mapping-coronavirus-responsibly/

Maine

Vt. | N.H.

Wash. | Idaho | Mont. | N.D. | Minn. | Ill. | Mich. | N.Y. | Mass.

Ore. | Nev. | Wyo. | S.D. | Iowa | Ind. | Ohio | Pa. | N.J. | Conn. | R.I.

Calif. | Utah | Colo. | Neb. | Mo. | Ky. | W.Va. | Va. | Md. | Del.

Ariz. | N.M. | Kan. | Ark. | Tenn. | N.C. | S.C. | D.C.

Okla. | La. | Miss. | Ala. | Ga.

Hawaii | Alaska | Texas | Fla.

Map | Cartogram

Lead **Win** Flip
Democrat

Lead **Win** Flip
Republican

https://www.bbc.co.uk/news/election-us-2020-54783016

3 electoral college votes

29 electoral college votes

IL

PA

NY

CA

DC

TX

FL

HI    AK

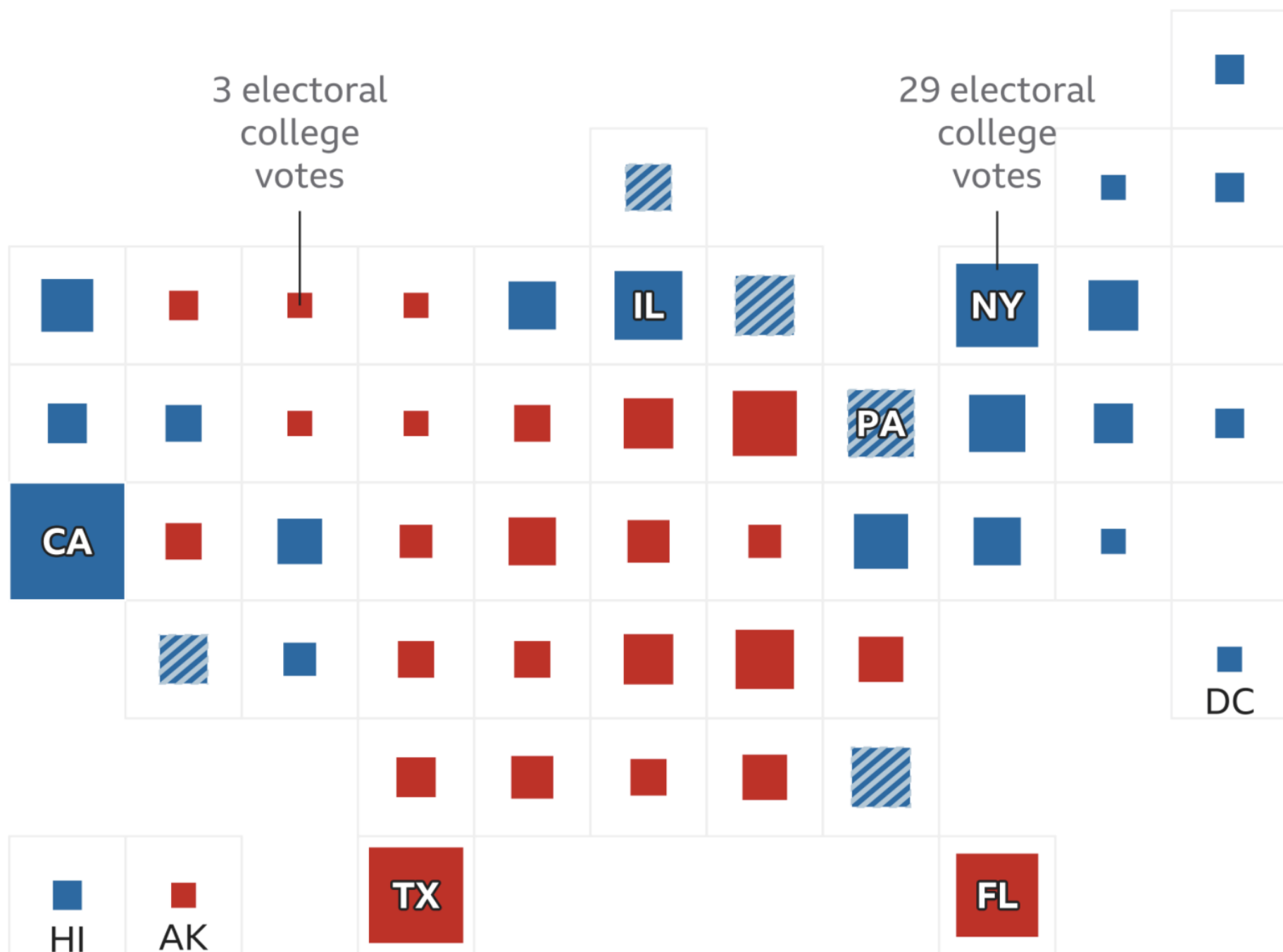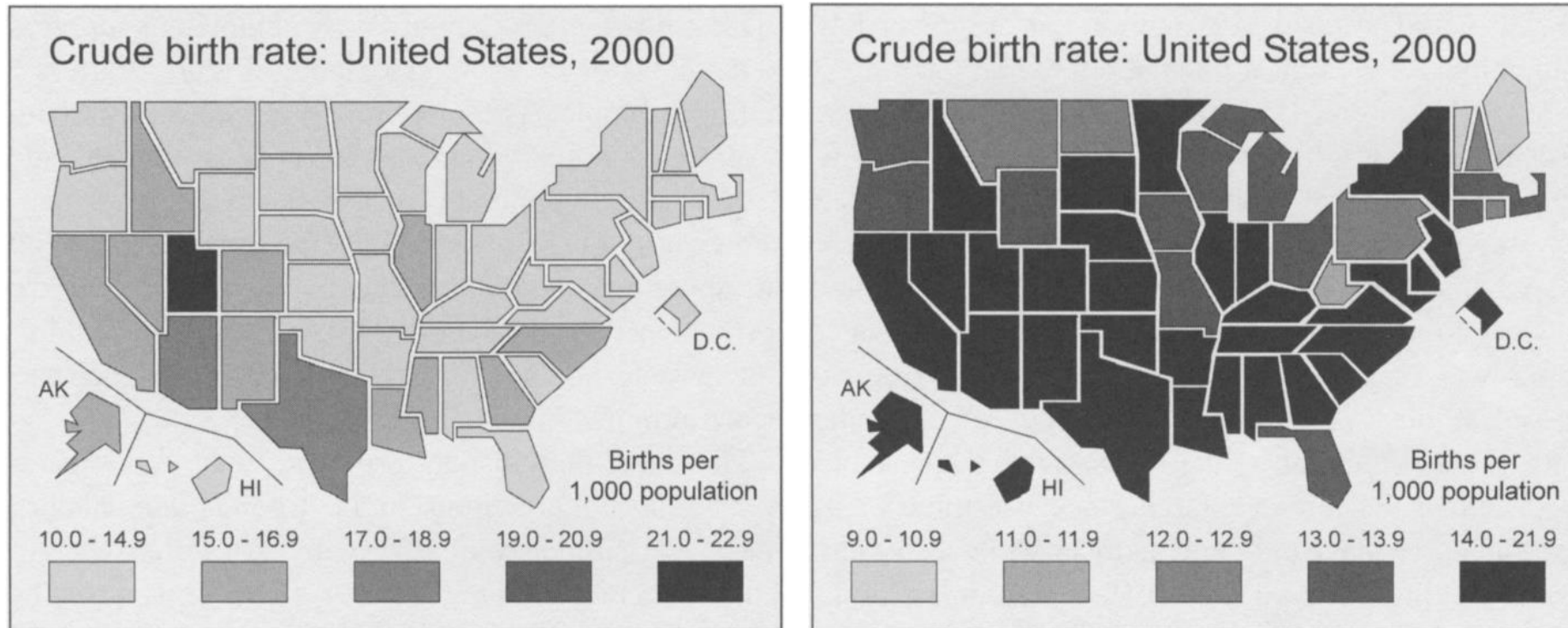# The first thing people do is looking at a visualisation
## (not the legend)



Choropleth maps: "darker-means-more metaphor"

# Now some DIY …