# Data Visualization

FOUNDATIONS (2)

Tea Tušar, Data Science and Scientific Computing

# The three principles of good visualization design

2

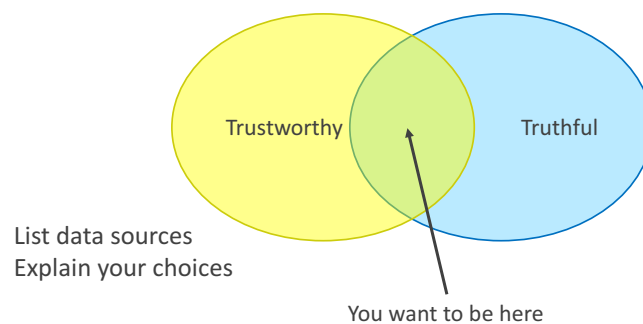# Good visualization design is

1. Trustworthy

2. Accessible

3. Elegant

A. Kirk. *Data Visualization*, SAGE Publications, 2016.　　　　　　　　　　　3
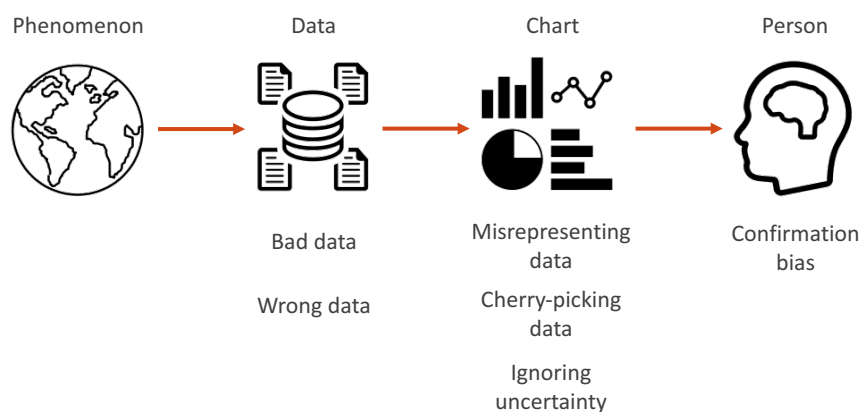
# Trustworthiness

Trust ≠ truth



Trustworthy　　　　　Truthful

List data sources
Explain your choices

You want to be here

4

# Trustworthiness

Lying with visualization is easy

*Intentionally and unintentionally*

5

# How charts lie?

| Phenomenon | Data | Chart | Person |
|---|---|---|---|

Bad data

Wrong data

Misrepresenting data

Cherry-picking data

Ignoring uncertainty
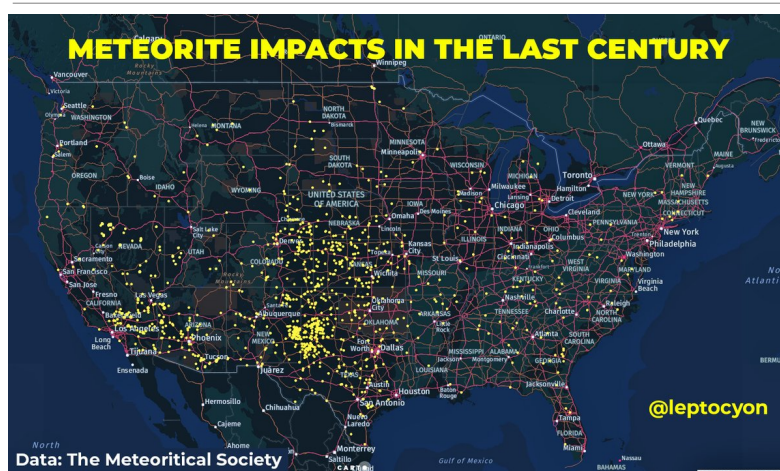
Confirmation bias

6

# Bad data

Garbage in, garbage out

- Unrepresentative data
  - Polls on unrepresentative populations
  - Measurements on unrepresentative samples
  - Too much missing data
- Biased data
  - Question framing in polls
  - Choice of measures

This is a problem when it is not made clear and the data is used for analyses that are suitable for more 'regular' data

7

# Unrepresentative data



Data: The Meteoritical Society

8

# Unrepresentative data

## Abraham Wald and the Missing Bullet Holes

Armour planes so that they don't get shot by enemy fighters. Armour is heavy, so use it only where is really needed.
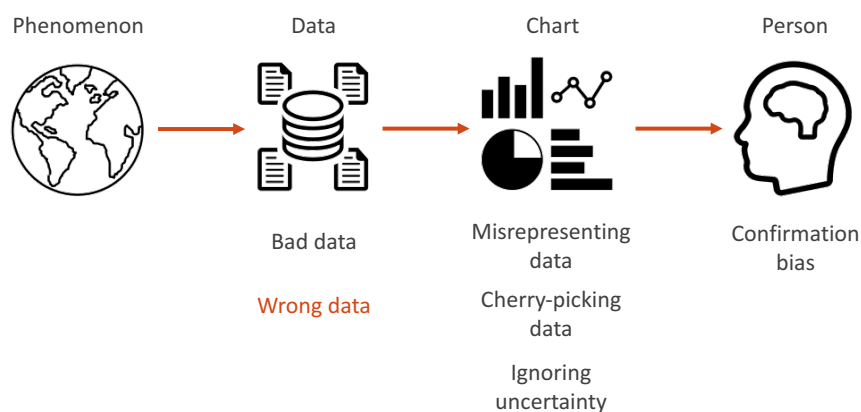
| Section of plane | Bullet holes per square foot |
|---|---|
| Engine | 1.11 |
| Fuselage | 1.73 |
| Fuel system | 1.55 |
| Rest of the plane | 1.8 |

https://medium.com/@penguinpress/an-excerpt-from-how-not-to-be-wrong-by-jordan-ellenberg-664e708cfc3d
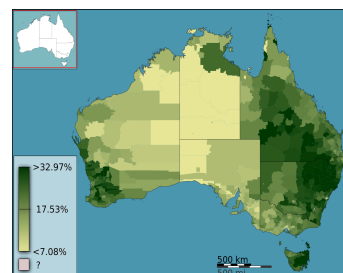
9

# How charts lie?

| Phenomenon | Data | Chart | Person |
|---|---|---|---|
| | Bad data | Misrepresenting data | Confirmation bias |
| | Wrong data | Cherry-picking data | |
| | | Ignoring uncertainty | |

10

# Wrong data

Comparisons using
- Non-comparable data
- Absolute instead of cumulative data (and vice versa)
- Absolute instead of relative data (in charts and choropleth maps)

The fraction of Australians that identified as Anglican at the 2011 census
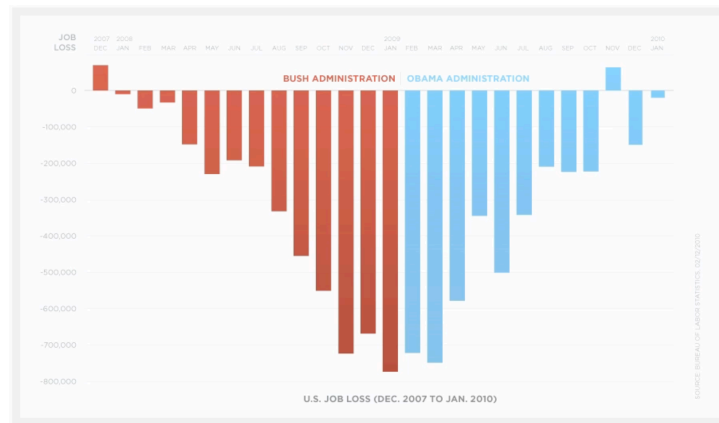
# Non-comparable data used in comparisons

Two issues
- Prices not adjusted for purchasing power
- Different sources of data

The data source specifically warns against using this data for comparison

# Absolute instead of cumulative data

# Absolute instead of cumulative data

# Absolute instead of cumulative data

# Cumulative instead of absolute data

# Cumulative instead of absolute data

17

# Absolute data in comparisons

18

# Absolute data in choropleth maps

# How charts lie?

| Phenomenon | Data | Chart | Person |
|---|---|---|---|



| | Bad data | Misrepresenting data | Confirmation bias |
| | Wrong data | Cherry-picking data | |
| | | Ignoring uncertainty | |

# Misrepresenting data

Ignoring conventions
- Placement of dependent and independent variables
- Inverted y axis
- Unequal intervals
- Pie charts that do not add up to 100%

Abusing scales
- Bar charts with truncated axis
- Aspect ratio bias
- Dual axes
- Improper scaling of areas (and pictograms)

Unnecessary 3-D

Improper categorization

Oversimplifying

21

# Inverted y axis

22

## Unequal intervals

## Unequal intervals

# Unequal intervals

# Over 100% pie chart

## Bar chart with truncated axis

## Aspect ratio bias

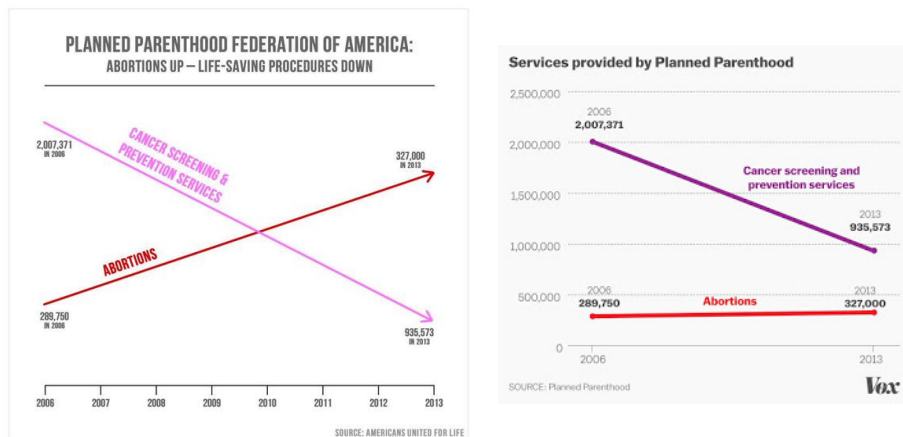

Banking to 45 Degrees

## Dual axes

## Improper scaling of areas (and pictograms)

# Improper scaling of areas (and pictograms)



Even worse, if the elements are 3-D

# Unnecessary 3-D

16

# Improper categorization

Figure 1. Prevalence of maternal smoking at any time during pregnancy, by state:
United States, 2016

WA
MT
ND
MN
OR
ID
WY
SD
WI
MI
NY
VT
ME
NH
MA
RI
CT
NV
UT
NE
IA
PA
NJ
DE
MD
DC
CA
CO
KS
MO
IL
IN
OH
WV
VA
AZ
NM
OK
AR
KY
TN
NC
SC
MS
AL
GA
TX
LA
FL
AK
HI

U.S. prevalence: 7.2%

4.9
2.1
7.7
4.9

NOTE: Access data table for Figure 1 at: https://www.cdc.gov/nchs/data/databriefs/db305_table.pdf#1.
SOURCE: NCHS National Vital Statistics System, Natality.

https://www.cdc.gov/nchs/data/databriefs/db305.pdf

35

# Improper categorization

States sorted by percentage

https://www.cdc.gov/nchs/data/databriefs/db305_table.pdf#1

36

18

# Improper categorization

# Oversimplifying

Clarify, not simplify!

*To clarify, add detail.*

Edward Tufte

# Box plot vs. violin plot

## Box (and whisker) plot

## Violin plot
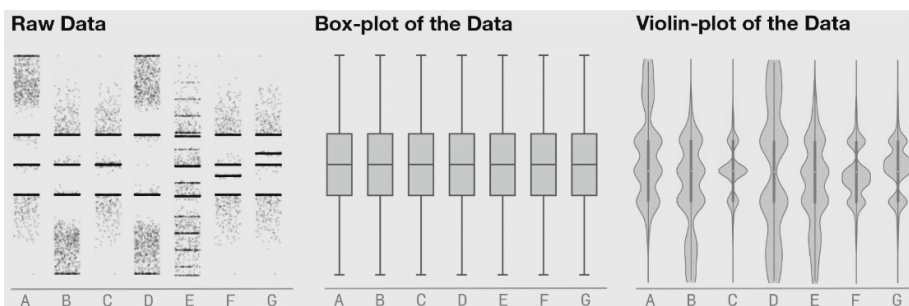
39

# Oversimplifying

**Raw Data** | **Box-plot of the Data** | **Violin-plot of the Data**

40

# How charts lie?

| Phenomenon | Data | Chart | Person |
|---|---|---|---|



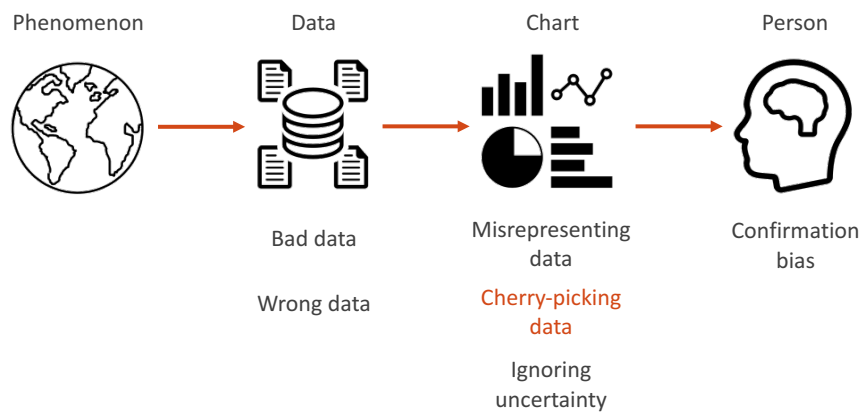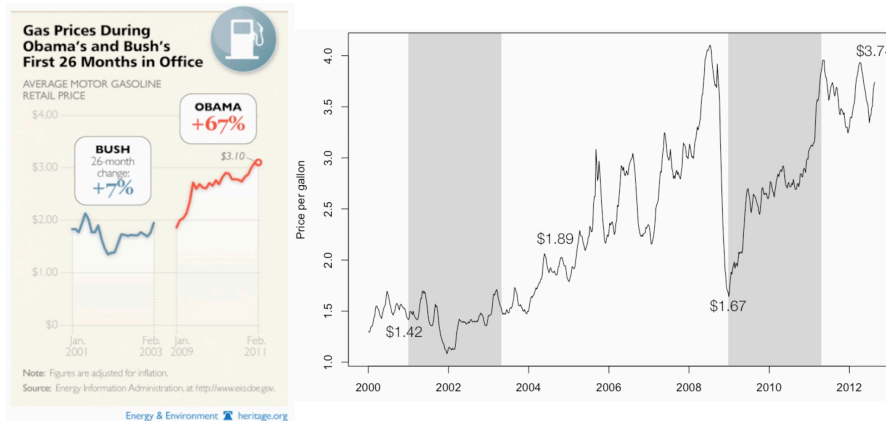| | Bad data | Misrepresenting data | Confirmation bias |
|---|---|---|---|
| | Wrong data | Cherry-picking data | |
| | | Ignoring uncertainty | |

41

# Cherry-picking data

A chart shows as much as it hides, so think about what might be missing

o Hiding (unfavorable) data
o Concealing existing patterns
o Suggesting patterns that are not there

Correlation ≠ causation

42

# Hiding (unfavorable) data

# Concealing existing patterns

Concealing existing patterns

https://blog.datawrapper.de/weekly47-cpi-dollars-for-college/ 45



Suggesting patterns that are not there

https://news.nationalgeographic.com/2015/06/150619-data-points-five-ways-to-lie-with-charts/ 46

# How charts lie?

| Phenomenon | Data | Chart | Person |
|---|---|---|---|



Bad data

Wrong data

Misrepresenting data

Cherry-picking data
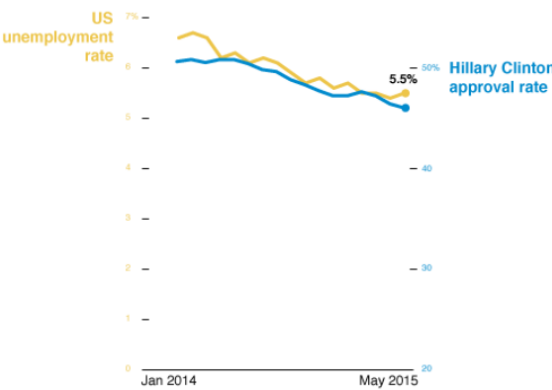
Ignoring uncertainty

Confirmation bias

49

---

# Ignoring uncertainty

o Misrepresenting uncertainty
o Concealing uncertainty

50

Data Visualization Foundations (2) header

# Misrepresenting uncertainty

The cone of uncertainty is widely misinterpreted

# Misrepresenting uncertainty

The cone of uncertainty is widely misinterpreted

# Misrepresenting uncertainty

The cone of uncertainty is widely misinterpreted

# Concealing uncertainty

# Concealing uncertainty

Directly age-standardised mortality from alcohol attributable conditions for men and women by borough in Surrey, rate per 100,000 people (2005/06).

59

# How charts lie?

| Phenomenon | Data | Chart | Person |
|---|---|---|---|



| | Bad data | Misrepresenting data | Confirmation bias |
| | Wrong data | Cherry-picking data | |
| | | Ignoring uncertainty | |

60

# Confirmation bias

Charts lie because we lie to ourselves – we see what we want to see

61

# Confirmation bias

# Confirmation bias

# Confirmation bias

# Confirmation bias



Surface on the
county-level map:
**Red: 80%**
**Blue: 20%**

Map by Kenneth Field
https://twitter.com/kennethfield/
status/970827334038237184

https://www.youtube.com/watch?v=Cd046xZhO_8&t=504s

65

# Confirmation bias



Surface on the
county-level map:
**Red: 80%**
**Blue: 20%**

**SHARE OF THE POPULAR VOTE IN THE 2016 PRESIDENTIAL ELECTION**

| | | |
|---|---|---|
| Donald Trump | **46.1%** | 62,984,825 votes |
| Hillary Clinton | **48.2%** | 65,853,516 votes |
| Other candidates | **5.7%** | |

**PERCENTAGE OF ELIGIBLE VOTERS**

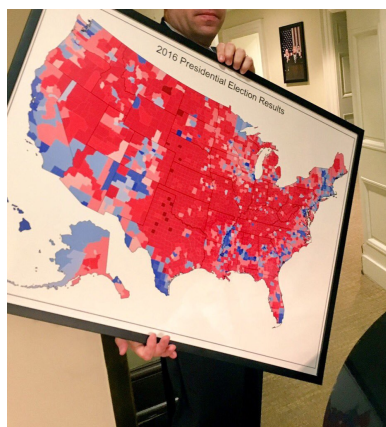| | |
|---|---|
| Didn't vote | **40.0%** |
| Voted for Donald Trump | **27.7%** |
| Voted for Hillary Clinton | **28.9%** |
| Voted for other candidates | **3.4%** |

https://www.youtube.com/watch?v=Cd046xZhO_8&t=504s

66

Confirmation bias

Surface on the county-level map:
**Red: 80%**
**Blue: 20%**

Bubble size is proportional to the number of votes received just by the candidate who won on each county

Confirmation bias

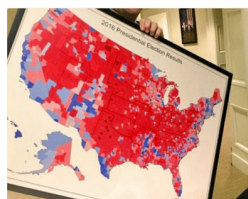SHARE OF THE POPULAR VOTE IN THE 2016 PRESIDENTIAL ELECTION

| | | |
|---|---|---|
| Donald Trump | 46.1% | 62,984,825 votes |
| Hillary Clinton | 48.2% | 65,853,516 votes |
| Other candidates | 5.7% | |

PERCENTAGE OF ELIGIBLE VOTERS

| | |
|---|---|
| Didn't vote | 40.0% |
| Voted for Donald Trump | 27.7% |
| Voted for Hillary Clinton | 28.9% |
| Voted for other candidates | 3.4% |

VOTES FOR DONALD TRUMP          VOTES FOR HILLARY CLINTON

Bubble size is proportional to the number of votes per county

# Confirmation bias

These are the numbers that truly matter in a U.S. Presidential Election

**ELECTORAL VOTES**   **TRUMP 304**   Other: 7   **CLINTON 227**

270

**WHO WON ON EACH STATE**   **STATE SIZE ADJUSTED BY ELECTORAL VOTES IT CONTRIBUTES TO THE ELECTION**



https://www.youtube.com/watch?v=Cd046xZhO_8&t=504s
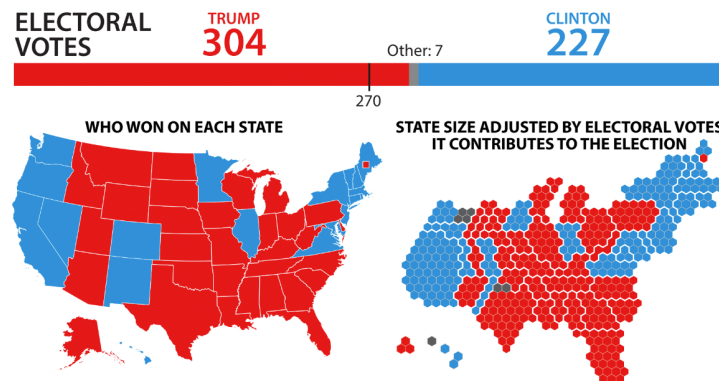
69

# To achieve trustworthiness

o List the source(s) of data
o Show representative and unbiased data (or clearly denote and explain why this is not the case)
o Compare only data that can be meaningfully compared
o Be mindful of the choice between absolute and cumulative values
o Use relative instead of absolute data in comparisons
o Follow conventions
o Do not abuse scales
o Do not use 3-D representations for non 3-D data
o Choose categories mindfully
o Do not oversimplify
o Present the entire relevant data
o Do not suggest patterns that are not there
o Show uncertainty
o Be wary of confirmation bias

70

# Trustworthiness

However… some rules can be bent (as long as you know what you are doing)

71