

Week 2: Graphics and Visualization

MATH-517 Statistical Computation and Visualization

Linda Mhalla

September 20th 2024

“The simple graph has brought more information to the data analyst’s mind than any other device.” – John W. Tukey

“The greatest value of a picture is when it forces us to notice what we never expected to see.” – John W. Tukey

One can think of graphics (and also models, for that matter) as a low-dimensional representation for data

Anscombe's Quartet

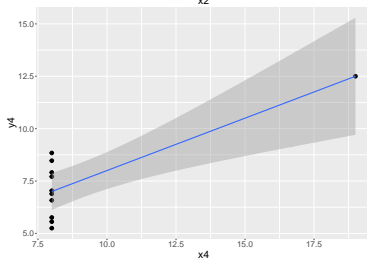
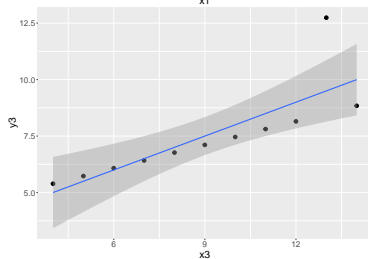
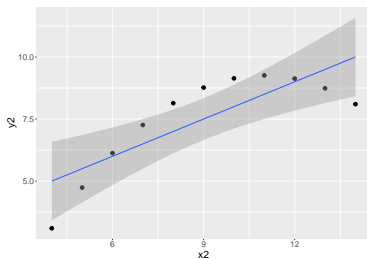
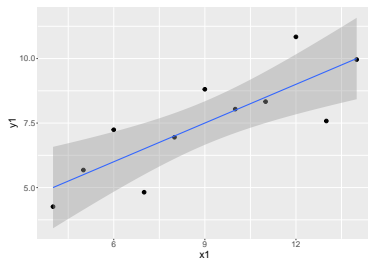
Four data sets with very similar descriptive statistics, each with

- one response variable y
- one regressor x

```
##      (Intercept)          x R-squared
## lm1      3.000091 0.5000909 0.6665425
## lm2      3.000909 0.5000000 0.6662420
## lm3      3.002455 0.4997273 0.6663240
## lm4      3.001727 0.4999091 0.6667073
```

⇒ how to fool the linear regression model ...

Anscombe's Quartet



⇒ qualitatively different structures

Datasaurus Dozen

A group of twelve datasets with almost identical summary statistics but are distinctively different when plotted!

Purposes of Visualization

The most common purposes of a visualization is

- data insight (exploratory graphs)
 - large data
 - detect important trends/patterns
 - find strange observations
- presentation (explanatory graphs)
 - result communication
 - decision making

⇒ a good choice of axes, axis limits, labels and symbols can facilitate substantially the extraction of information from the data

Datasets used for illustration are described in details in the [Lecture Notes 2](#)

What makes a good graph?

“... that which gives the viewer the greatest number of ideas in the shortest time with the least ink in the smallest space” – Edward Tufte

*“... graphical excellence requires telling the truth about the data.”
– Edward Tufte*

Tufte's design principles for effective data visualization

- Principle 1: maximize data-ink ratio
- Principle 2: minimize chart junk
- Principle 3: minimize lie factor

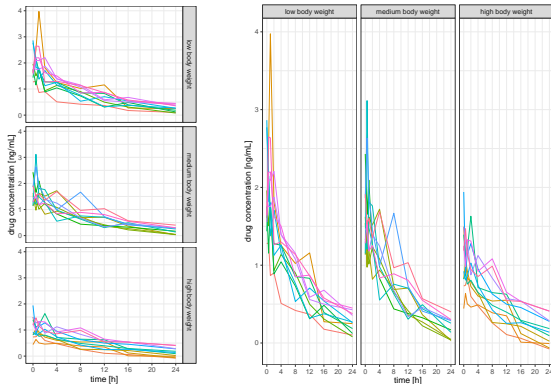
Section 1

Elements of Visualization

Elements of Visualization: Layout

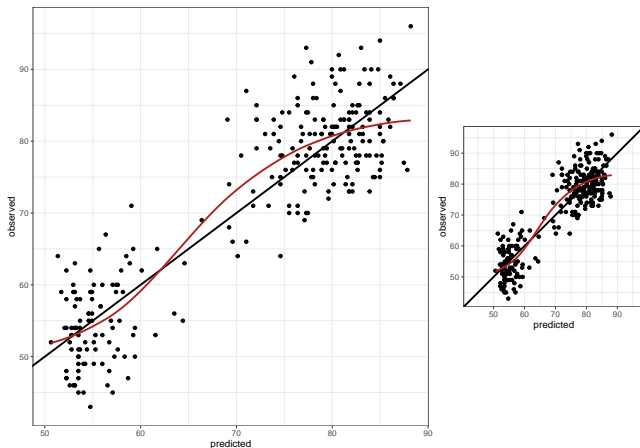
The arrangement of multiple panels is used for efficient comparison. To compare data on

- the x-axis, stack panels sharing a single x-axis
- the y-axis, use a single y-axis with all panels aligned horizontally



Elements of Visualization: Aspect Ratio

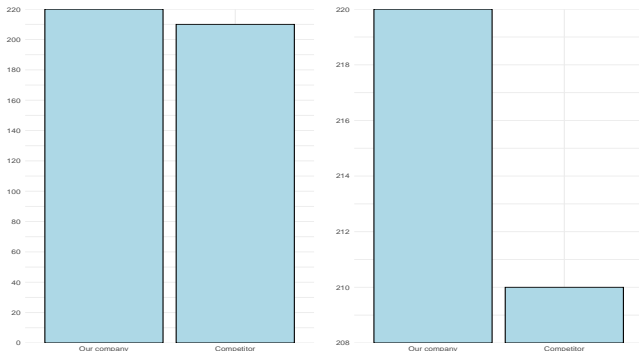
- If measurements on both axes reflect the same quantity (e.g., before vs after treatment, observed vs modelled), a square figure (1:1 aspect ratio) could avoid visual bias + make distances consistent and comparable



Elements of Visualization: Axes

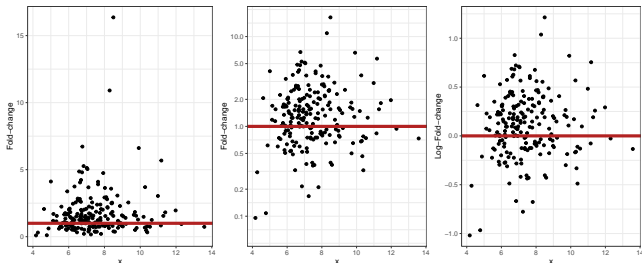
- If values are positive, axes should start at zero (unless there is a good reason for different choice) and not contain negative values

→ this usually applies to barplots used to highlight relative differences



Elements of Visualization: Axes

- If values are ratios or relative changes, axes should be logarithmic and symmetric around the point of no change (e.g., 0)



Elements of Visualization: Colours and Shapes

Chart features (ggplot's arguments) for points (and similarly for lines)

- `colour`: help identify different groups (do not use for pure decoration)

Tufte pointed out that because “they do have a natural visual hierarchy, varying shades of gray show varying quantities better than color”, and “the shades of gray provide an easily comprehended order to the data measures. This is the key”.

Elements of Visualization: Colours and Shapes

- colour
- shape: if data are ordered, use ordered symbols (number of vertices)

```
pch = _  
1 ○ 6 ▽ 11 ⌘ 16 ● 21 ◉  
2 △ 7 ☒ 12 ⊞ 17 ▲ 22 ▣  
3 + 8 * 13 ⊗ 18 ◆ 23 ◇  
4 × 9 ⊕ 14 ⊞ 19 ● 24 ▲  
5 ◇ 10 ⊕ 15 ■ 20 • 25 ▽
```

- size
- alpha (opacity/transparency): can be used to clarify plots with many points

These features can be used

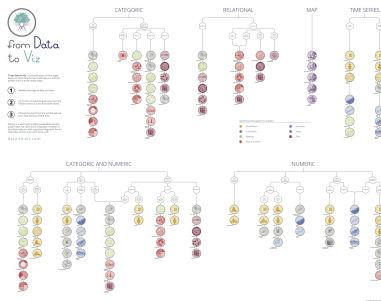
- to include additional information (or dimensions, i.e., to include additional variables) in a scatterplot
- to combat overplotting

Section 2

Visualization types

What type of chart?

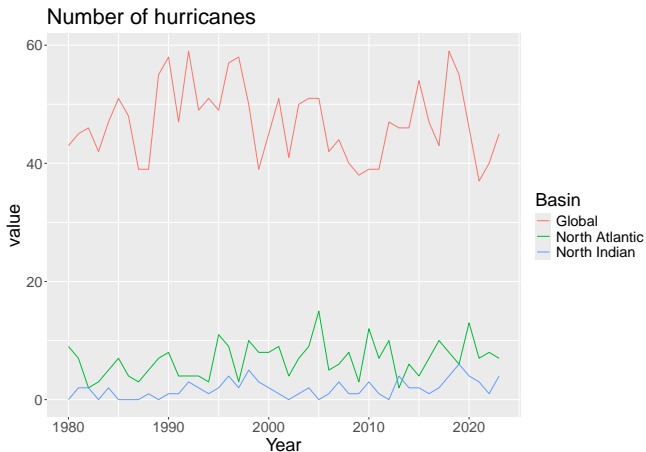
- Be clear about the intended purpose of a data visualization (compare groups, highlight correlation, show evolution over time, ...)
- Choose the right chart according to the type of data and purpose of visualization: [From Data to Viz](#)



- Further general guidance can be found in [Robbins' Creating Better Graphs](#)

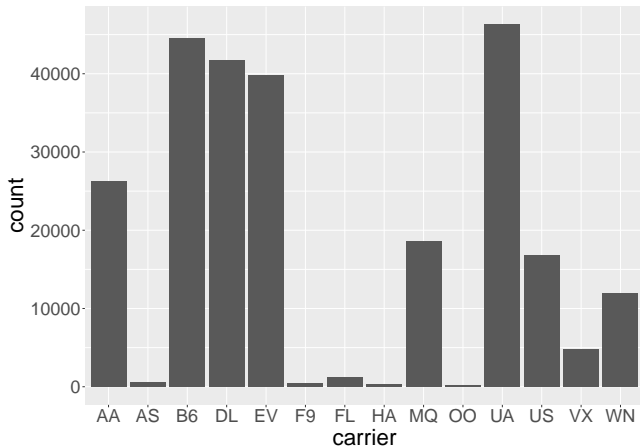
Line Plot

- there needs to be a linearly ordered variable, typically time
- if multiple groups with an inherent order, choose line types with an order (e.g., of thickness or dash density)



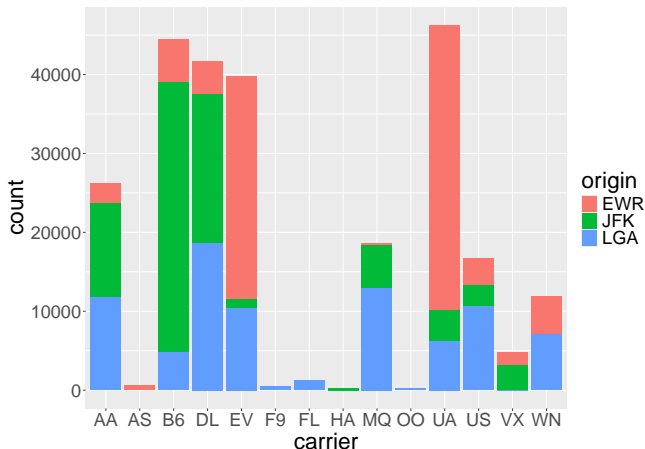
Bar Plot

- shows the relationship between a numeric and a categorical variable



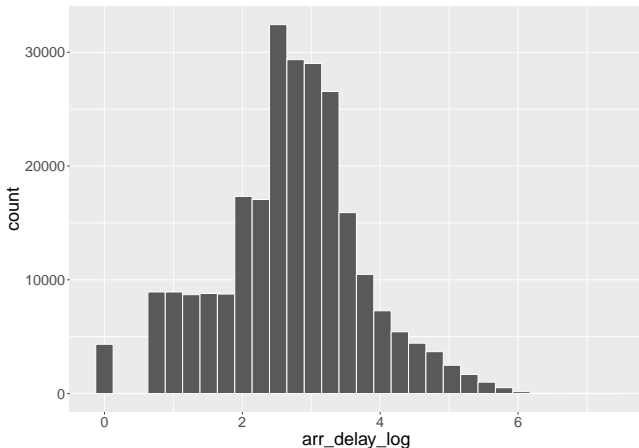
Bar Plot

- shows the relationship between a numeric and a categorical variable
- or displays values for several grouping levels



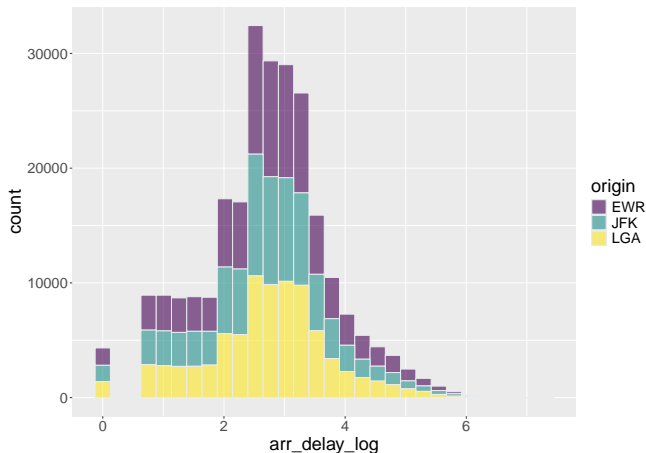
Histogram

- the variable is cut into several bins, and the number of observations per bin is represented by the height of the bar
- graphical representation of the distribution of a numeric variable
- allows to check the distribution for mistakes, outliers, or extremes



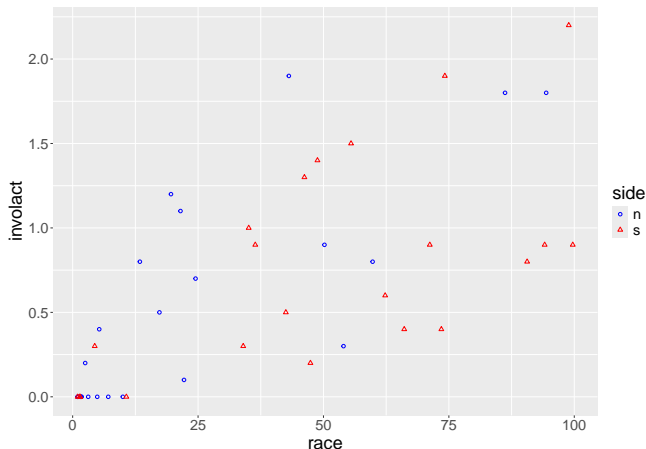
Histogram

- histogram allows to compare the distribution of several variables (not too many...)



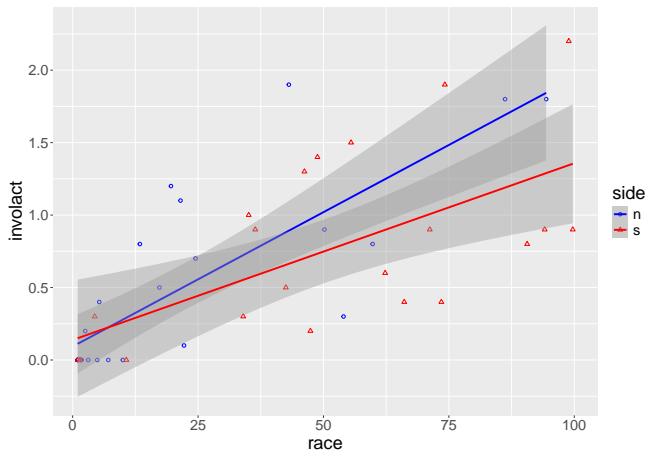
Scatterplot

- allows to study the relationship between two numeric variables
- subgroups can be added for more insights into hidden patterns

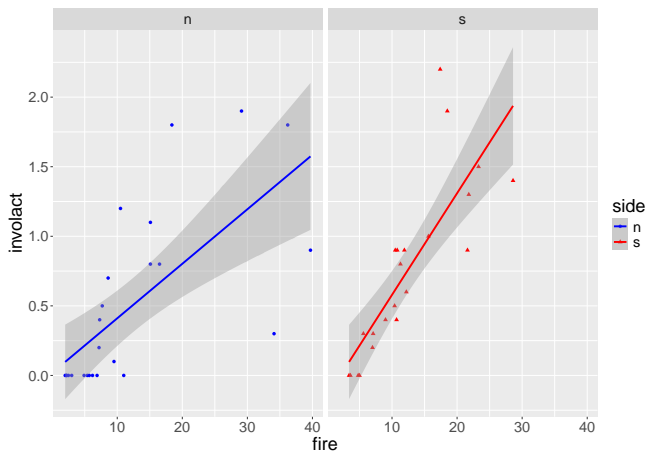


Scatterplot with Regression Lines

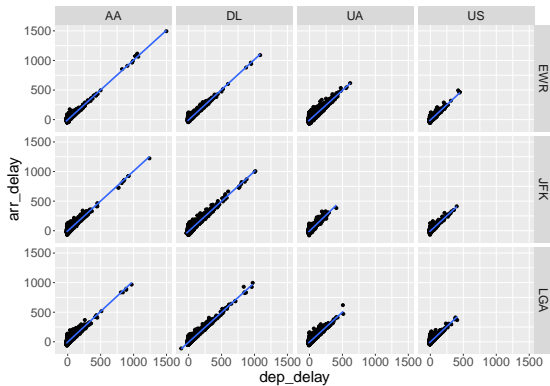
- different types of relationships can be detected by adding fitted regression curves



Plot by Factor

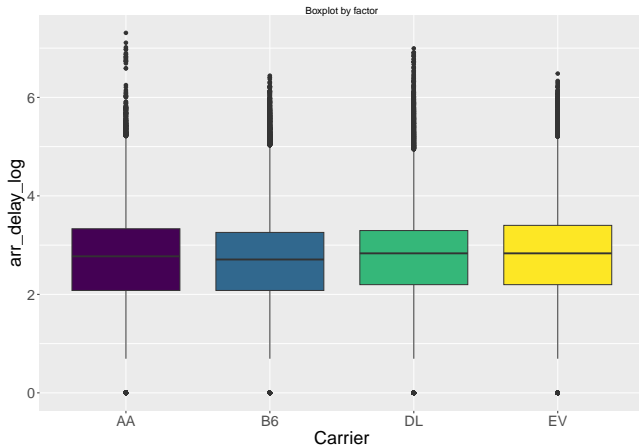


Plot by Two Factors

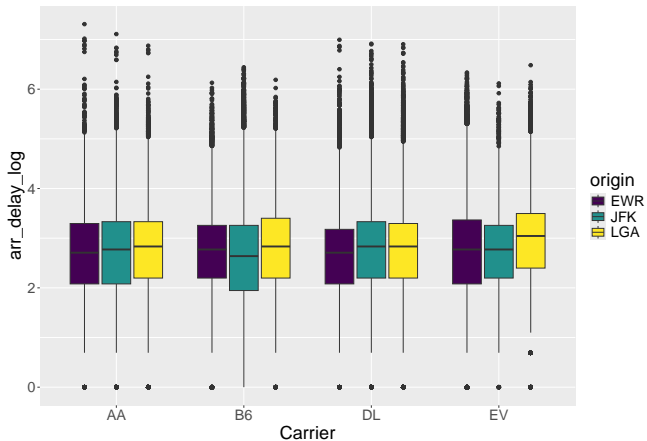


Boxplot

- summarizes the distribution of a numeric variable for several groups (presence of outliers, symmetry, dispersion, ...)

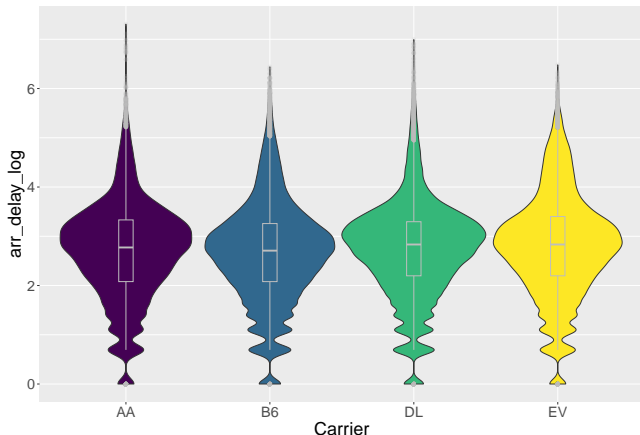


Boxplot by Two Factors



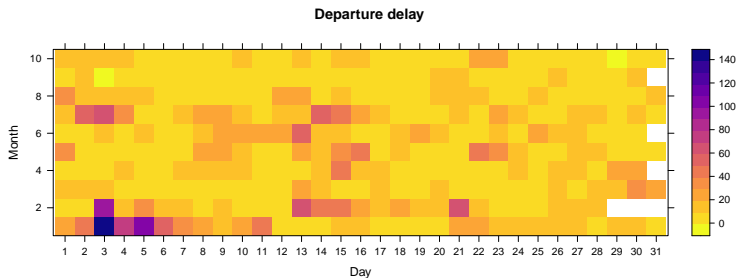
Improved version of a boxplot: The violin plot

- adds a representation of the distribution (a smooth version of the histogram)



Heatmaps

- represents a large matrix of data with colours reflecting their values (widely used for gene expression data)
- can be combined with dendrograms to display to show clustering, e.g., of genes



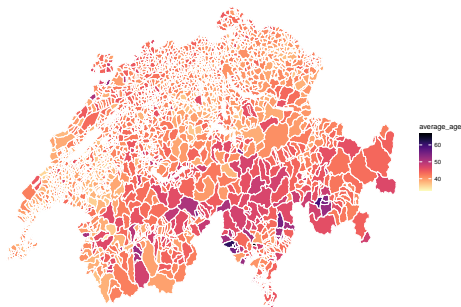
- could also be used for correlation matrices

Visualization of Spatial Data

Spatial data are complicated due to different

- data structures (vector or raster data)
- data sources
- data processing packages
- visualization packages

(Not so) short course about visualizing spatial data [here](#) (only if interested)

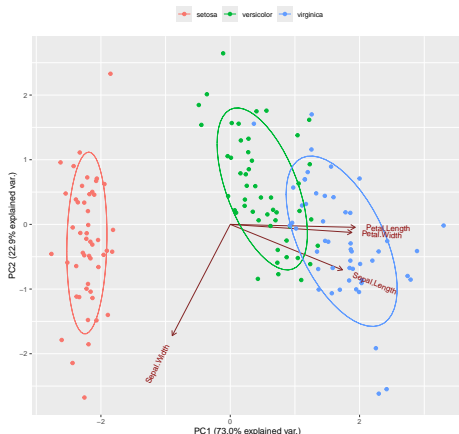


inspired by [this blogpost](#)

Visualization of High-dimensional Data

Exploration of high-dimensional data can be done using

- clustering methods: clusters are formed based on similarities between features (e.g., K -means, hierarchical clustering, ...)
- PCA: dimension reduction technique that preserves most of the data variability



Good Visualization Practices

- context: exploratory vs. explanatory analysis/graphics
 - exploratory ... helps you understand the patterns
 - explanatory ... designed to communicate your understanding
- provide context (in text **and** in caption)
- take the audience into account: yourself? readers of your paper? listeners of your talk? what is their background?
- ask yourself if the graph
 - answers a/your question
 - achieves your aim (of convincing?)

Good Visualization Practices

- gray scale often preferable
 - color-blindness (friendly palettes, e.g. [Coolors](#))
- axes (scale, gaps, etc.)
 - text of appropriate size
- publication-specific conditions
- seek simplicity and clarity, though you can be artistic!
 - sometimes bend the rules (responsibly and justifiably). As [Nigel Holmes](#) stated: “As long as the artist understands that the primary function is to convey statistics and respects that duty, then you can have fun (or be serious) with the image; that is, the form in which these statistics appear.”
 - unusual complex graphs are [more memorable](#)
 - But, remember that visualization must be **accurate**, **easy to comprehend**, and **appropriate to the context**

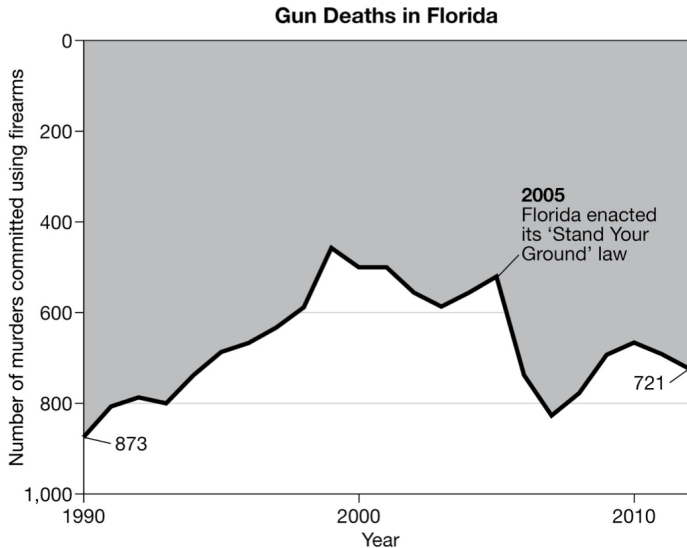
Find inspiration in [The R Graph Gallery](#).

Beware when exporting graphics

Section 3

Bad Visualization Practices

Reverted Axis

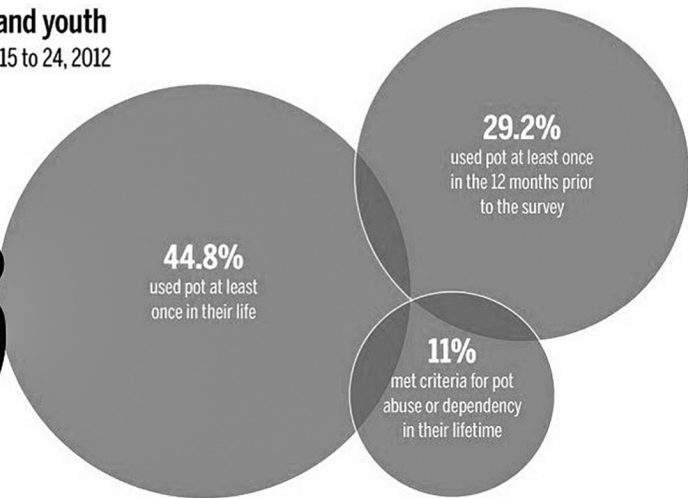


Source: Florida Department of Law Enforcement

False Venn's Diagrams

Marijuana and youth

Canadians age 15 to 24, 2012



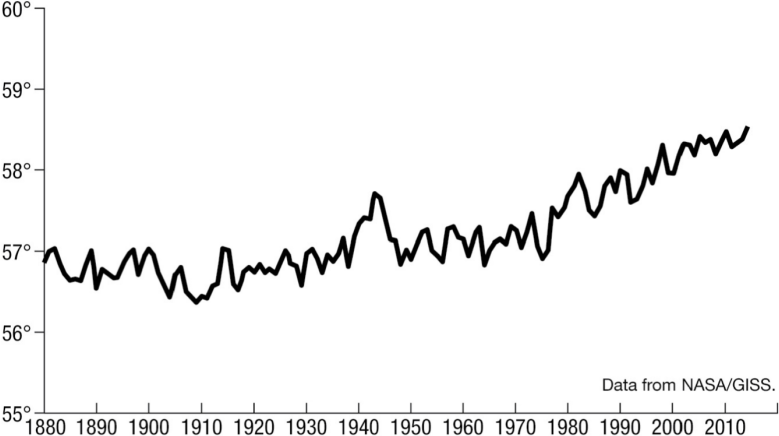
SOURCE: STATISTICS CANADA

Axis Starting at Zero

Average Annual Global Temperature in Fahrenheit, 1880–2019

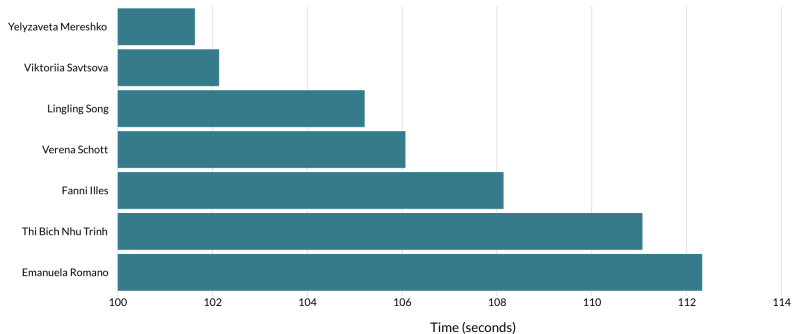


Average Global Temperature by Year

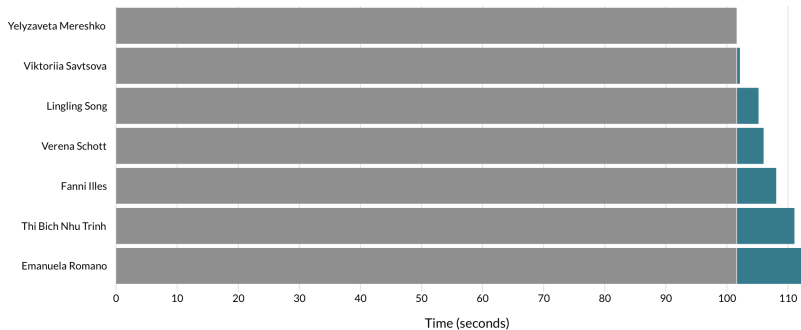


Over-exaggeration/ Missing Baseline

Women's 100m breaststroke SB5, Rio 2016



Women's 100m breaststroke SB5, Rio 2016



To cut or not to cut?

“In general, in a time-series, use a baseline that shows the data not the zero point” – Edward Tufte

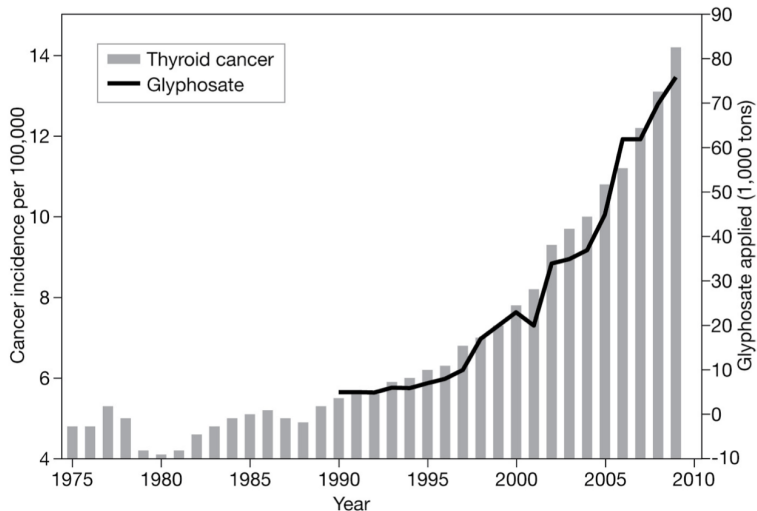
“don’t spend a lot of empty vertical space trying to reach down to the zero point at the cost of hiding what is going on in the data line itself” – Edward Tufte

So what?

- Barplot: With this kind of chart there is consensus: your Y-axis should start at 0
- Line plot: Here however there is no consensus, even if in general you don’t have to start at 0

Double Axes

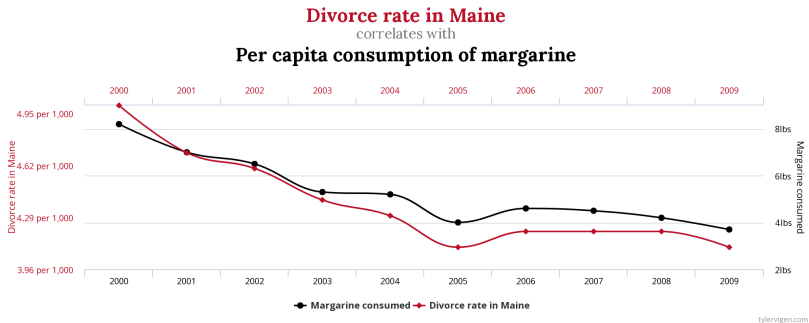
Thyroid Cancer Incidence Rate



- This is actually quite good, but double axes are usually problematic

False Causation

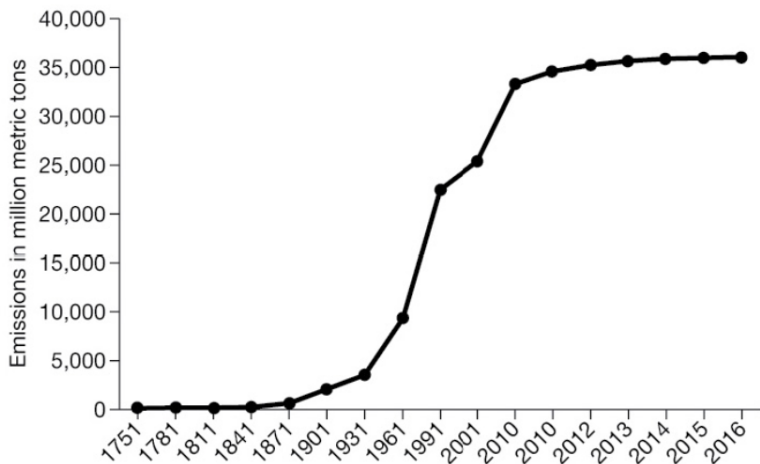
Correlation does not imply causation



Data sources: National Vital Statistics Reports and U.S. Department of Agriculture

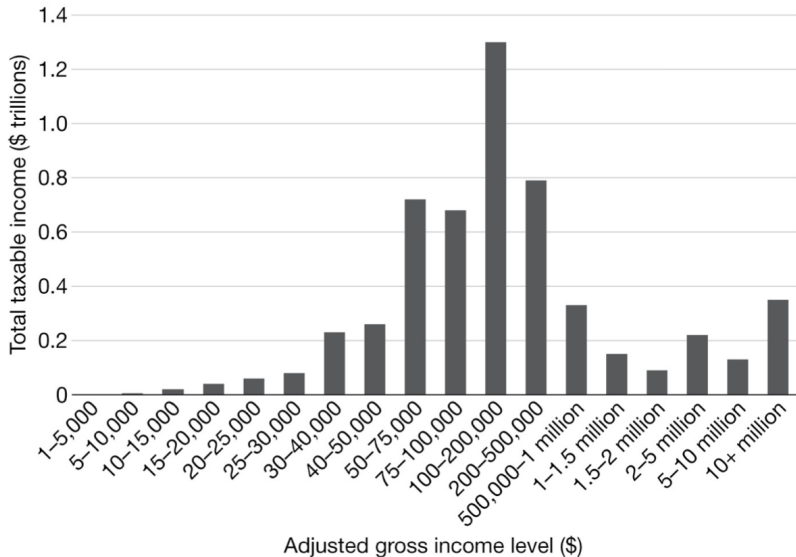
Tweaking Axis

Carbon Dioxide Emissions from Global Fossil Fuel Combustion and Industrial Processes, 1751–2016



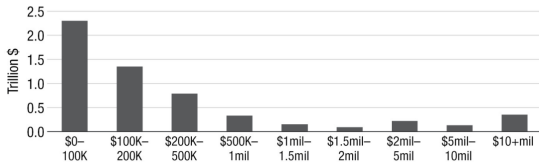
Binning

Total Taxable Income in 2008

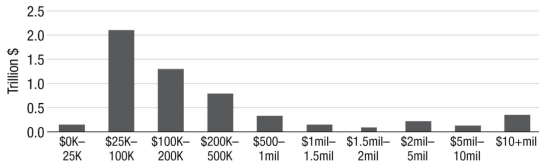


Different Kinds of Binning

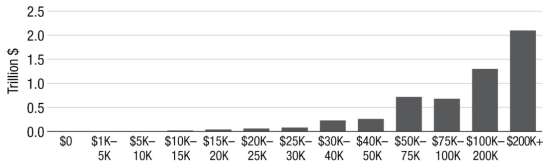
Tax the Poor!



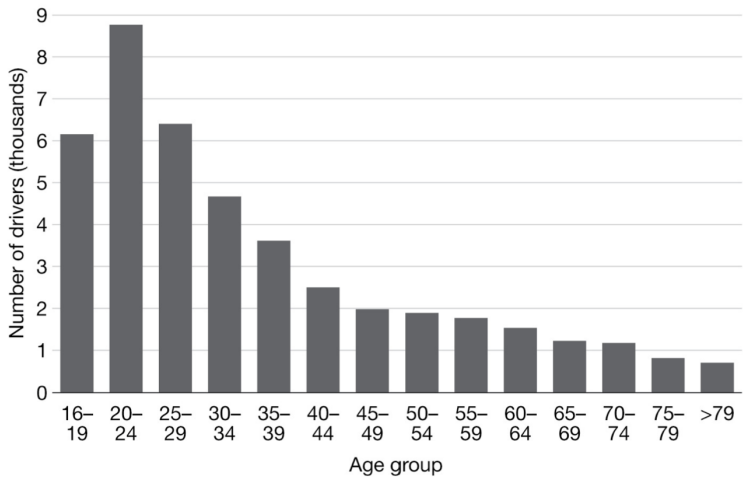
Tax the Middle Class!

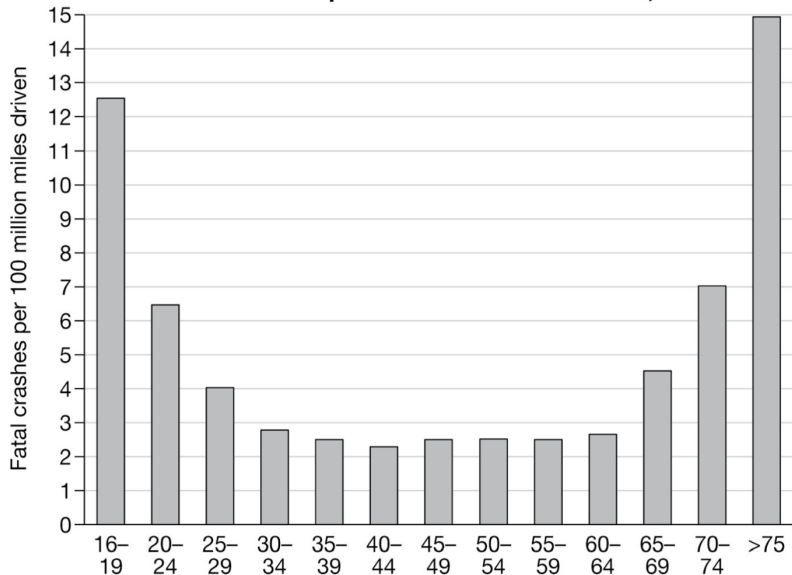


Tax the Wealthy!



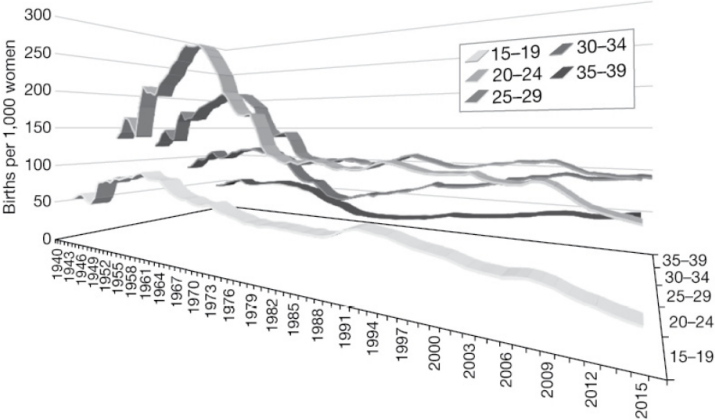
Number of Drivers in Fatal Crashes, 1988



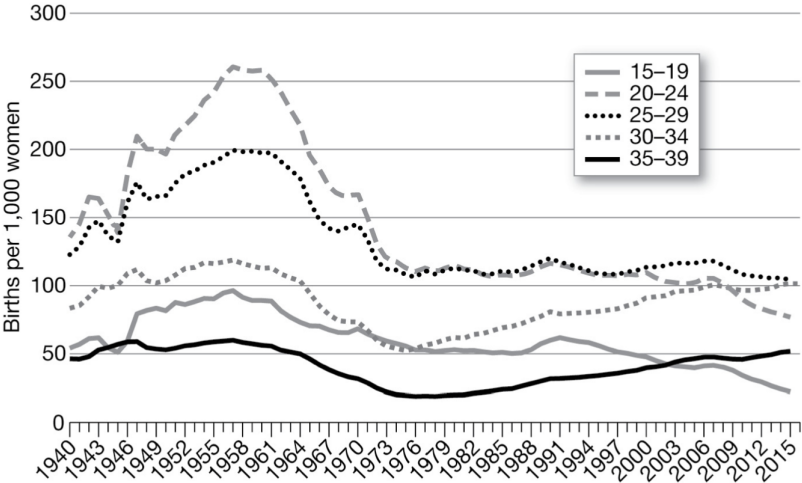
Fatal Crashes per 100 Million Miles Driven, 1988

Useless 3D

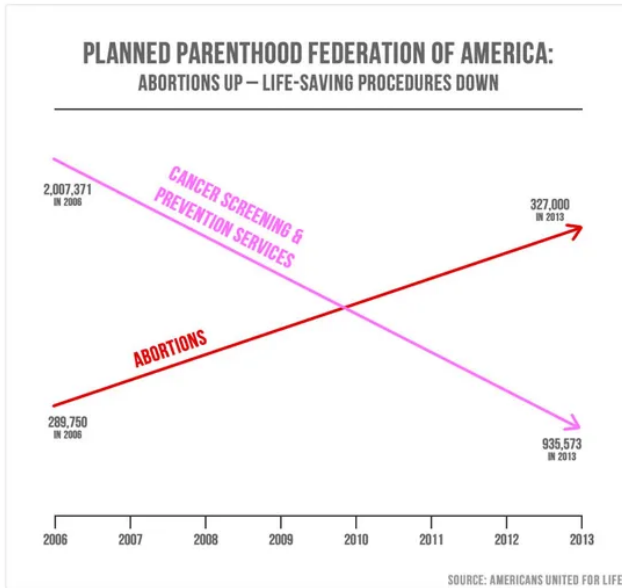
Female Birth Rates by Age, U.S.A.



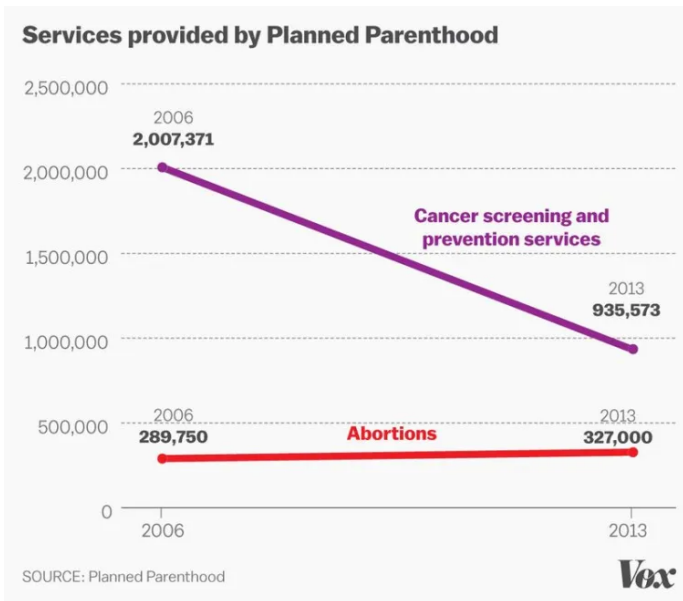
Female Birth Rates by Age, U.S.A.



Missing Axis & Misguidance (1)



Missing Axis & Misguidance (2)



Data Visualization Caveats

A collection of caveats in [From Data to Viz](#)

Assignment

Small project [20%]. Deadline on Week 5.

The goal of this project is *data exploration*. Details can be found on the dedicated [course page](#).

Some links to open data can be found [here](#).

References

- Krause, Andreas, Nicola Rennie and Brian Tarran (2023) [Best Practices for Data Visualisation](#)
- Poldrack (2019) [Statistical Thinking for the 21st Century](#)
- [JASA Ethical Guidelines for Statistical Practice](#)
- Gelman (2018) [Ethics in statistical practice and communication](#)
- Wickham & Golemund (2017) [R for Data Science](#)