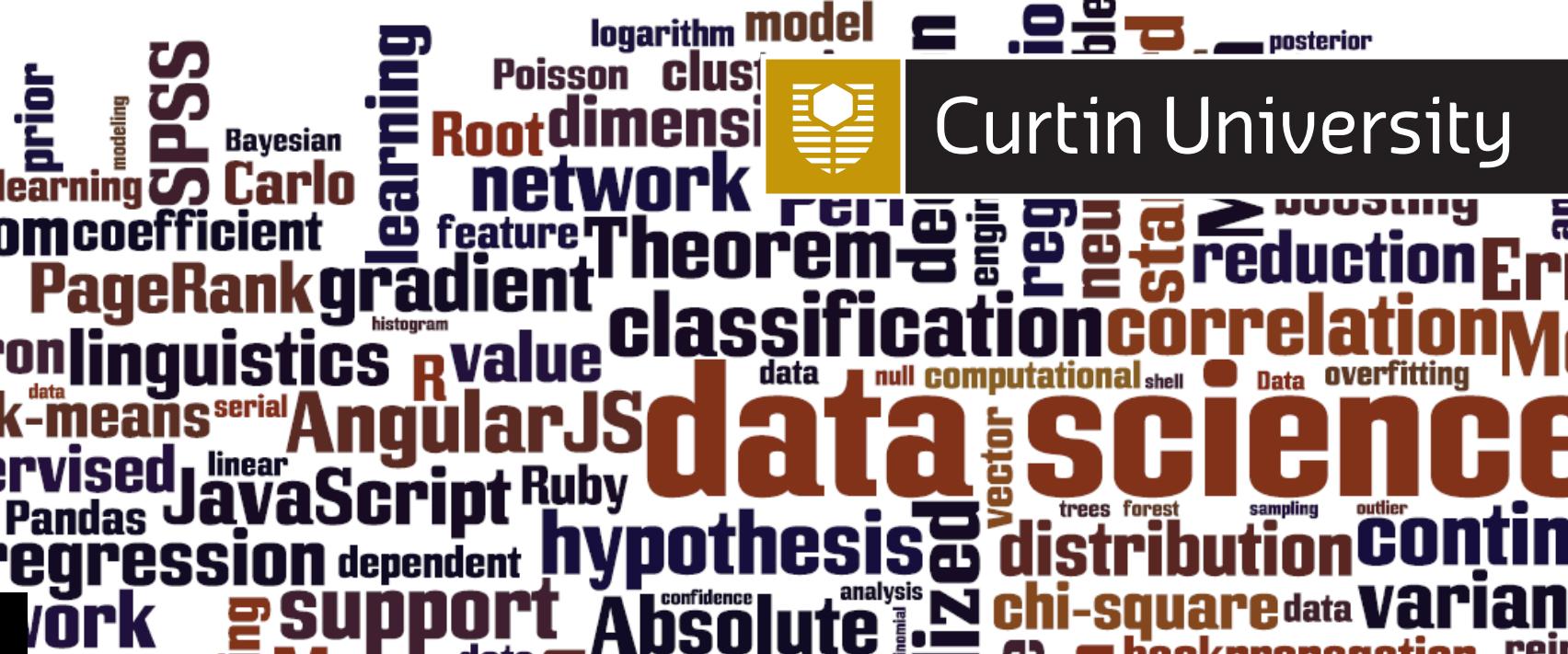


Lecture 3



# STAT1003 Introduction to Data Science



# Test 1

**Date/Time/Location:** Teaching Week 6, during your computer lab. Please come to the time slot in which you're registered.

**Duration:** ~90 minutes

**Format:** You will download and then work in an R Markdown file that will contain the questions (just as you do during computer labs). You will also be required to download some data, either from BB or the web. Please make sure you know where to save these files so you have access to them via RStudio. If you don't have an I:\ drive (it should begin with your student number) please see CITS. Please also practice knitting the R Markdown file to a Word file. At the end of the test, you'll be submitting both to BB.

**Content:** Short-answer questions based on lectures, lecture notes, and the readings thus far (Weeks 1 - 3); Questions that will require you to manipulate data, answer questions about it, and produce plots; There may be one question (~10% of overall mark) that will challenge/extend your existing knowledge and capabilities in R

**Aids:** Help in R/Rstudio; Lecture notes, your own notes, readings, and computer labs & their solutions; NO searching the web or going to external websites during the test

**Practice Test:** Will be uploaded to Teaching Week 3 Unit Materials by Friday, solution next week.

**Drop-in:** Times TBA- Wednesday afternoons, room 314.462

The e-book *Data Wrangling with R* is quite a useful resource for learning about manipulating data.

Week 3/3

STAT1003 Lecture 3



Curtin University

# Outline

- Visualization
  - Role of visualization
  - Visualization is not new!
  - Elements of visualization
- Additional examples

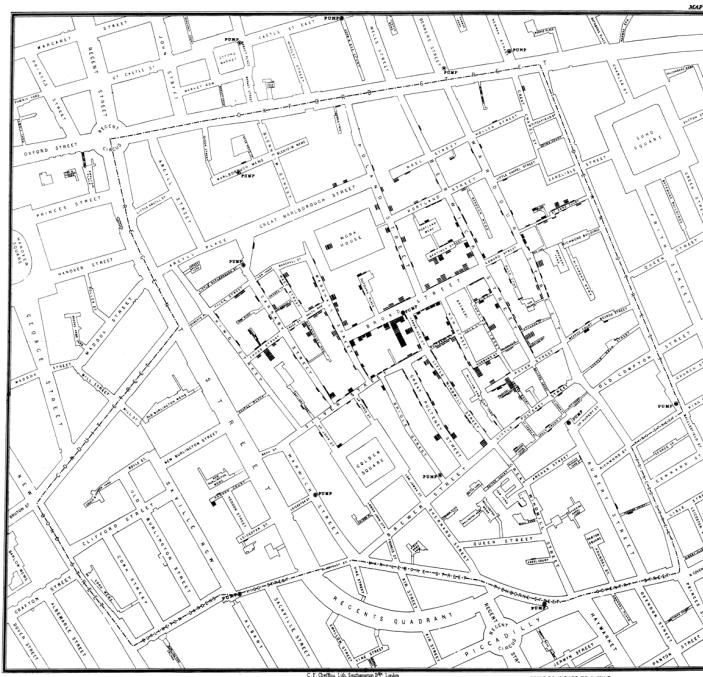
# Visualization and graphics

- Visualization starts with the question, “What do I want to know about my data?” but it is an iterative process
- Visualization belongs in every stage of the investigation
  - Exploratory data analysis
    - Getting to know the data
    - Uncovering structure that cannot be seen in tables or numerical summaries
  - Modelling and analysis
    - Understanding results
  - Presenting results
    - Visualization is an essential part of communicating a message or story

# Visualization and graphics

- What type of variables/data are in the dataset?
  - Categorical or quantitative
  - Ordered or unordered categories?
  - Discrete or continuous?
  - Time series, spatial, spatio-temporal?
    - Equally or unequally spaced?
  - Other types of data?
- Graphics and visualization are not new!

# John Snow and the Broad Street pump (1854)

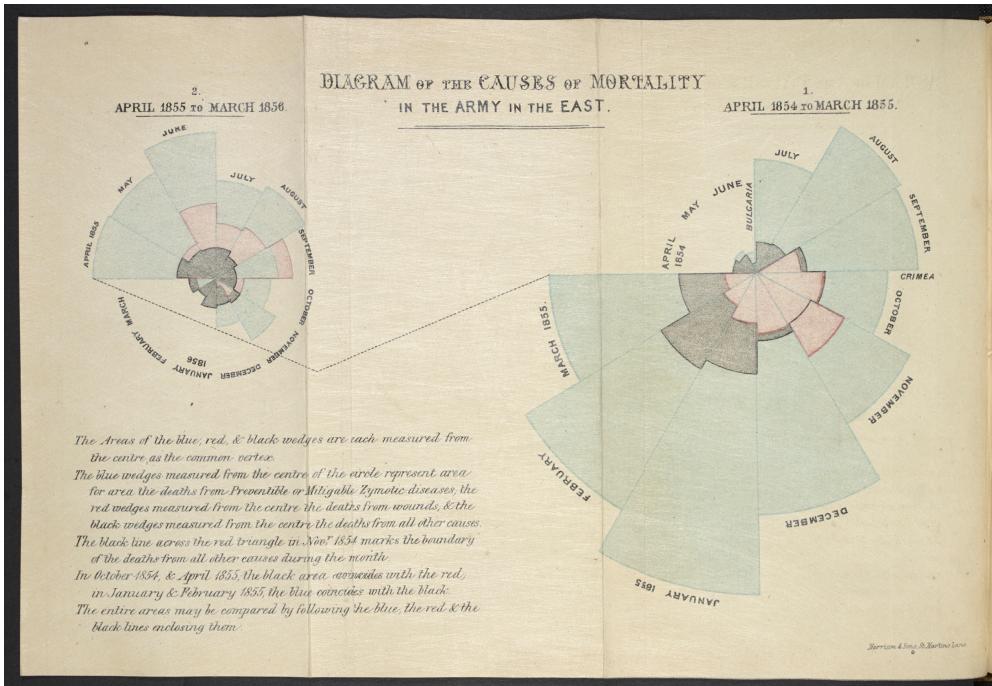


Week 3/3

STAT1003 Lecture 3



# Florence Nightingale's rose (1858)

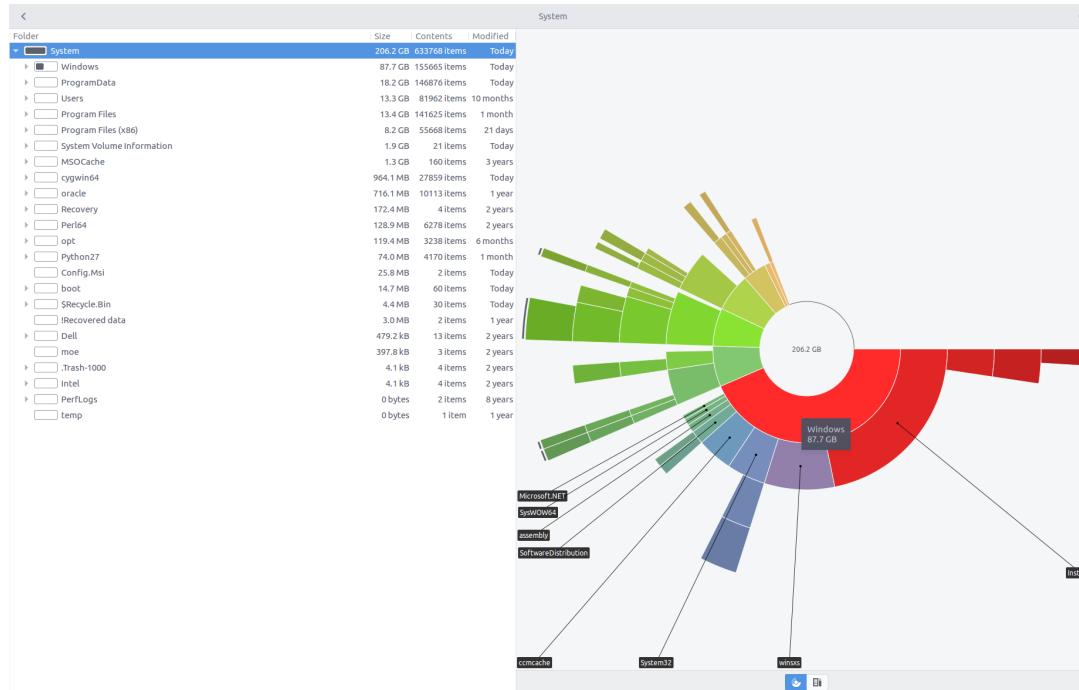


Week 3/3

STAT1003 Lecture 3



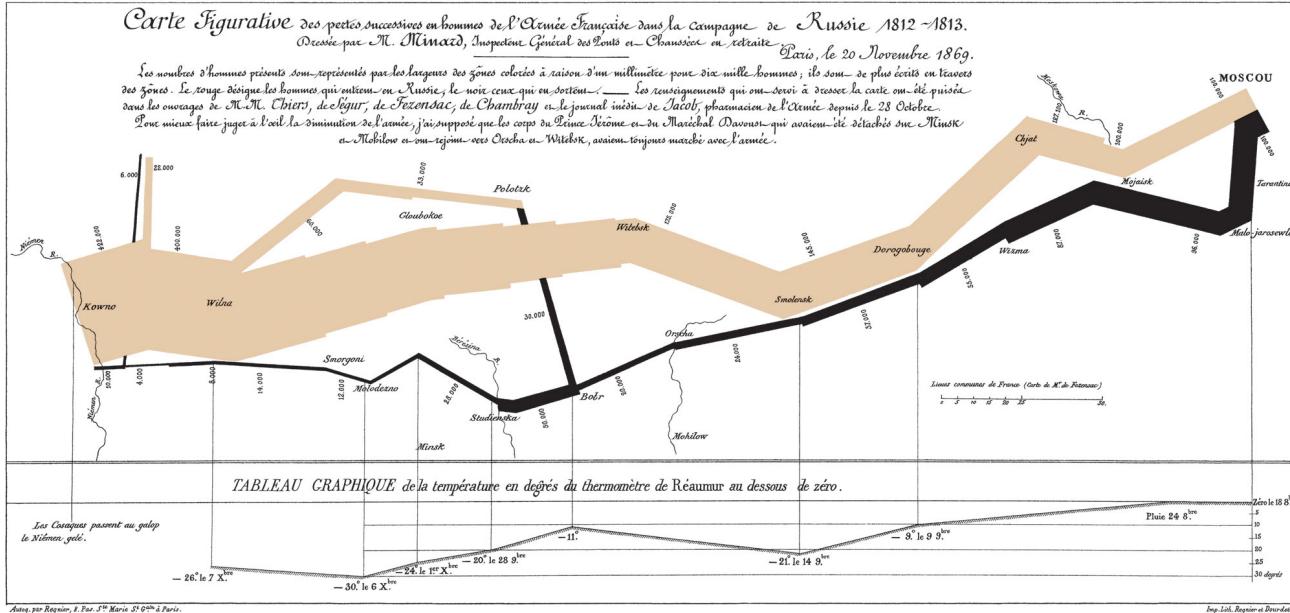
# Contemporary rose diagram



Week 3/3

STAT1003 Lecture 3

# Charles Minard (1869) and Napoleon's army



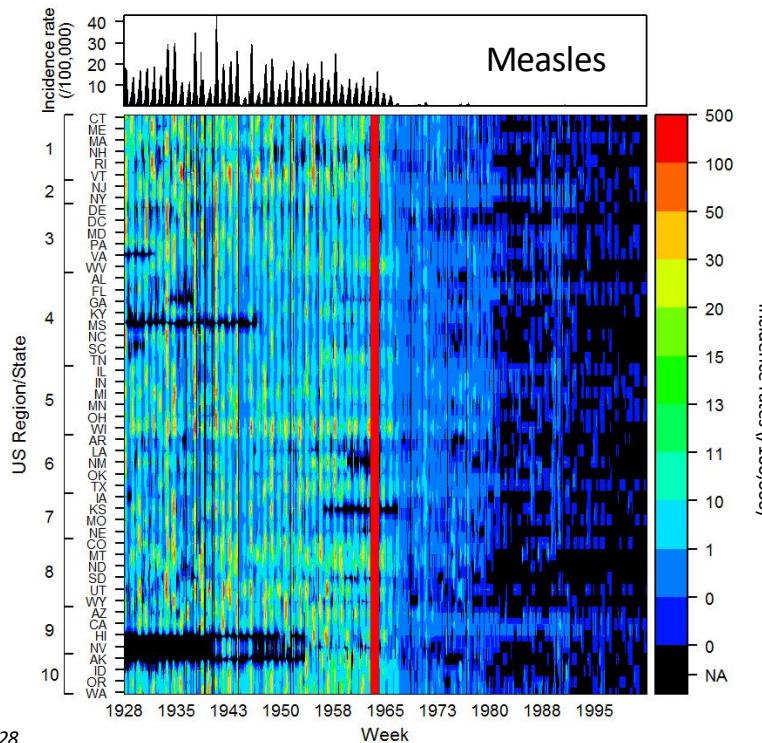
Week 3/3

STAT1003 Lecture 3



Curtin University

# A contemporary example: Decline of measles (I)



Van Panhuis WG, et.al., NEJM 2013 Nov 28

Week 3/3

STAT1003 Lecture 3

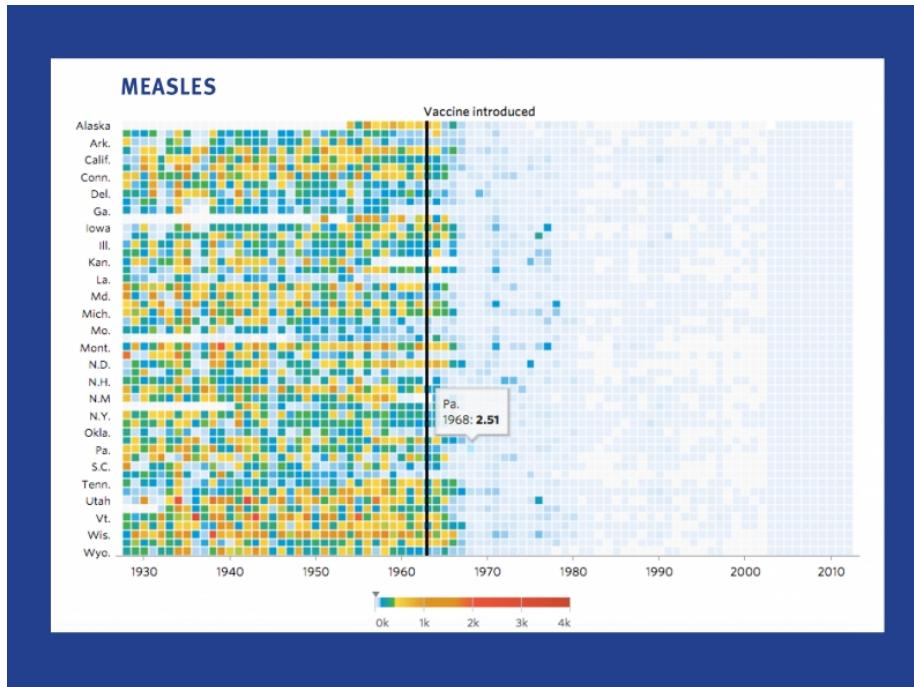


# Decline of measles (1)

- What is (are) the main message(s)?
- What is (are) the secondary message(s)?
- What elements are used?
- How is the graphic organized?
  - Context matters
- What comparisons are being visualized?

# Example – Decline of measles (II)

[http://www.pittmed.health.pitt.edu/sites/default/files/u48/13\\_Measles\\_B\\_big.jpg](http://www.pittmed.health.pitt.edu/sites/default/files/u48/13_Measles_B_big.jpg)



Week 3/3

STAT1003 Lecture 3



# Elements of visualization

- Scale
- Conditioning
- Perception – colour & length
- Transformations
- Context
- Smoothing and other large data considerations

*Deborah Nolan (2017)*

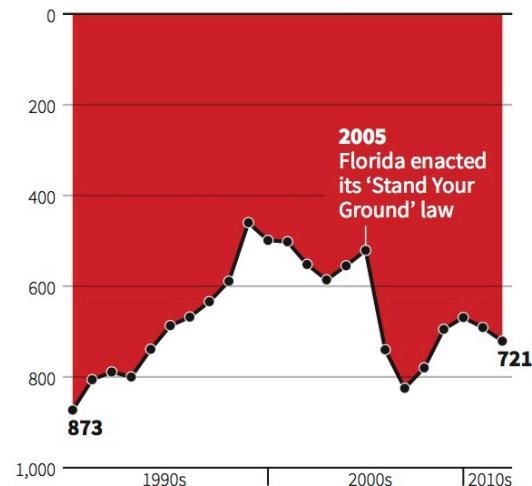
# Scale

- In any plots, choose axis limits to fill the plotting region
- If necessary, make multiple plots of different regions to focus on important features that may not be visible in a single plot
- Is the plot constructed in a way that matches viewers' expectations?
  - Do the axes increase in an expected way?
- Don't confuse the viewer!
  - Scales changes mid-axis
  - Two different scales on the same axis

# What's misleading about this?

## Gun deaths in Florida

Number of murders committed using firearms



Source: Florida Department of Law Enforcement

C. Chan 16/02/2014

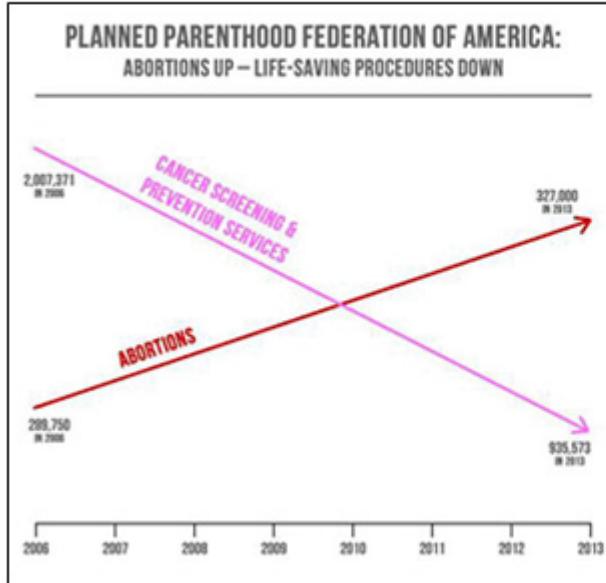
REUTERS

Week 3/3

STAT1003 Lecture 3



# Or this?



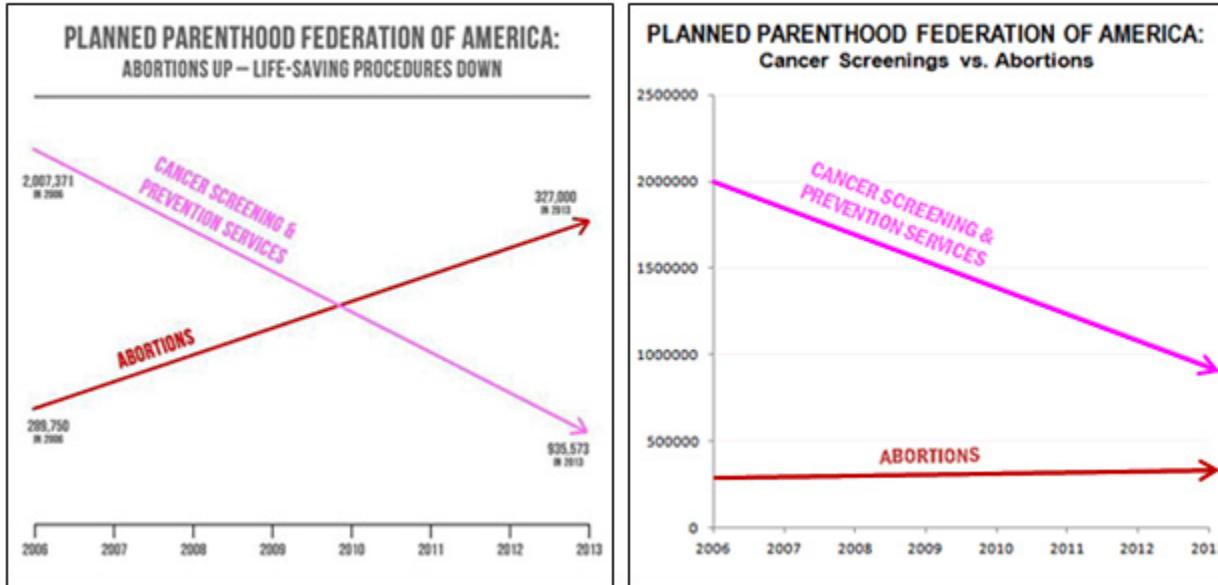
Week 3/3

STAT1003 Lecture 3



Curtin University

# Here's what's misleading



Week 3/3

STAT1003 Lecture 3

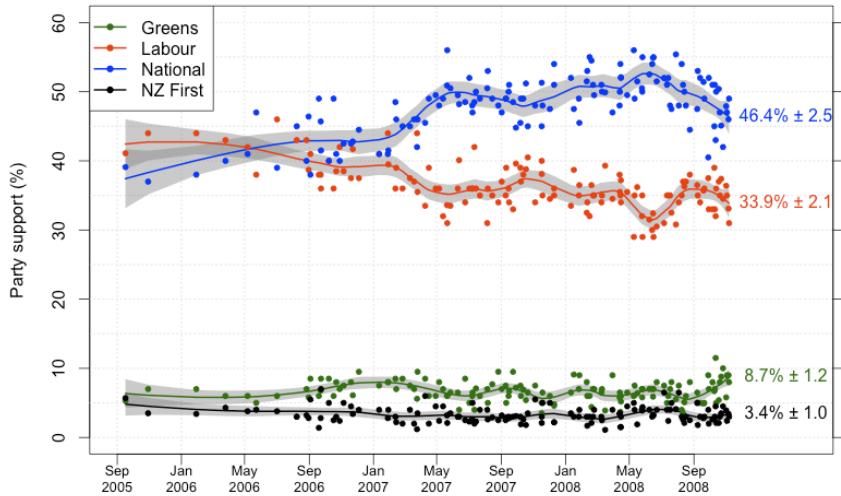


Curtin University

# Conditioning

- In any large and complex dataset, there will likely be different relationships among variables in different subgroups of the data
  - Construct side-by-side scatterplots, histograms, boxplots, etc. but keep scales the same to make comparisons easier
  - If representing on the same plot, use different colours and symbols to represent different groups
  - Use scatterplot smoothers to guide the viewer's eye
- In the measles example, some grouping according to region was carried out

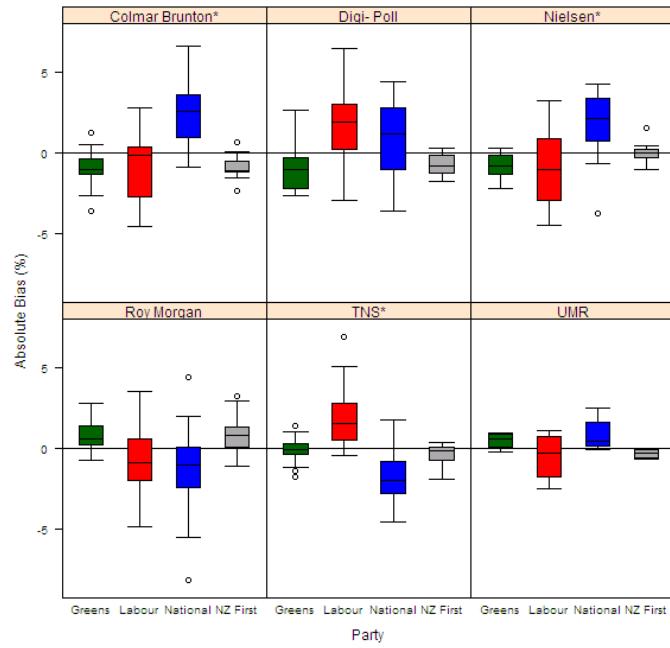
# Examples of conditioning



[Mark Payne](#)

Week 3/3

STAT1003 Lecture 3

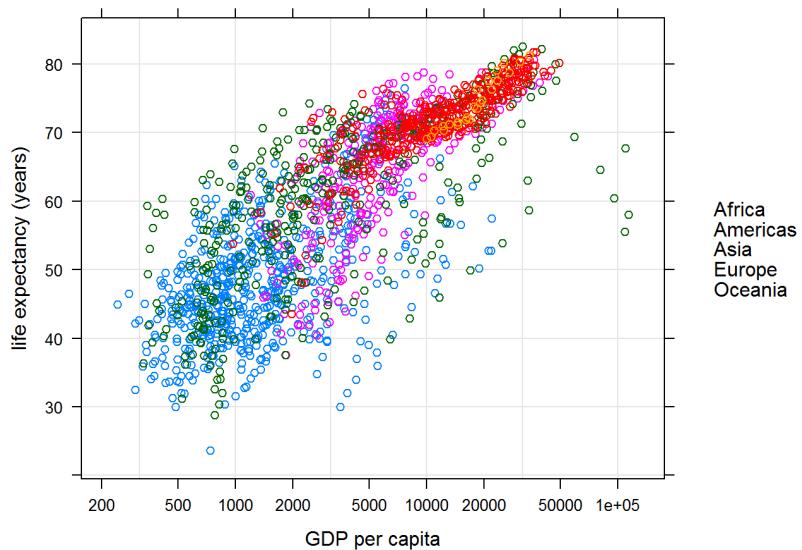


Companies whose name is marked with an asterisk (\*) show significantly different biases between the Labour and National parties at the 90% level (95% one tailed test).

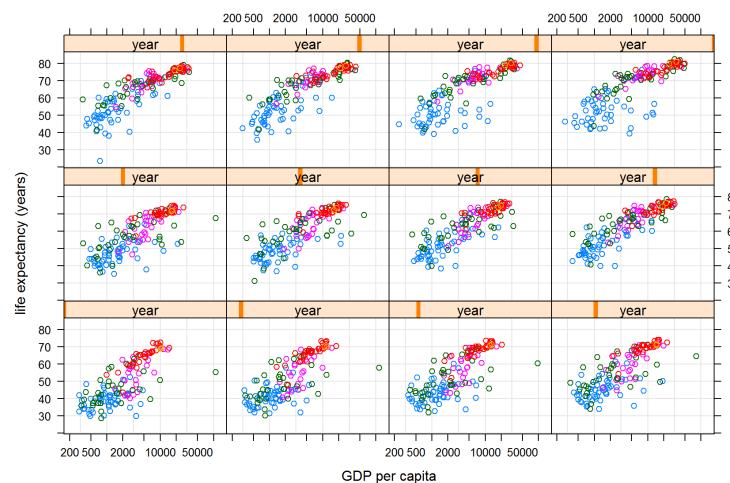


Curtin University

# Examples of conditioning



Africa  
Americas  
Asia  
Europe  
Oceania



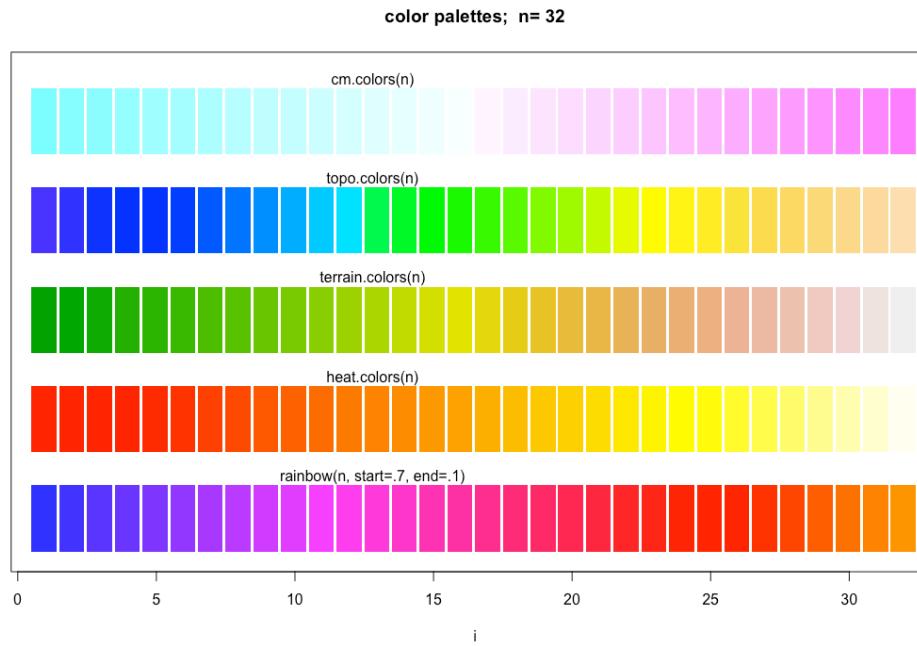
Africa  
Americas  
Asia  
Europe  
Oceania



# Perception – color

- Colour plots are much more common, but being able to choose a palette of colours is not a straightforward task
- Saturated (bright) colours tend to be difficult to look at for long periods of time
- Keep in mind that some people are colour-blind
- Use existing colour palettes rather than creating your own
- Decide whether to use a sequential colour scheme or one where low and high values should be highlighted?

# Color palettes in R



Week 3/3

STAT1003 Lecture 3



# Colour palettes in R

<https://www.nceas.ucsb.edu/~frazier/RSpatialGuides/colorPaletteCheatsheet.pdf>

**R color cheatsheet**

Finding a good color scheme for presenting data can be challenging. This color cheatsheet will help!

**H**ex to hexademical to represent colors

Hexadecimal is a base-16 number system used to describe color. Red, green, and blue are each represented by two characters (#rrggbb). Each character has 16 possible symbols: 0,1,2,3,4,5,6,7,8,9,A,B,C,D,E,F.

"00" can be interpreted as 0.0 and "FF" as 1.0 i.e., red = #FF0000, black = #000000, white = #FFFFFF

Two additional characters (with the same scale) can be added to the end to describe transparency (#rrggbbaa)

R has 657 built in color names Example: To see all lists of names: `colorNames`

These colors are displayed on P. 3.

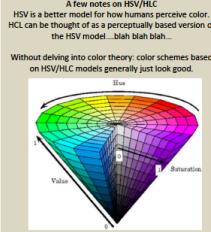
R translates various color models to hex, e.g.:

- RGB (red, green, blue): The default intensity scale in R ranges from 0-1; but another commonly used scale is 0-255. You can do this in R using `maxColorValue=255`. `alpha` is an optional argument for transparency, with the same intensity scale.
- rgb(r, g, b, maxColorValue=255, alpha=255)
- HSV (hue, saturation, value): values range from 0-1, with optional argument `alpha`, `h, v, s, alpha`.
- HCL (hue, chroma, luminance): hue from 0-360, 0 = red, 120 = green, blue = 240, etc. Range of chroma and luminance depend on hue and each other `hcl(c, l, alpha)`

A few notes on HSV/HLC

HSV is a better model for how humans perceive color. HCL can be thought of as a perceptually based version of the HSV model... blah blah blah...

Without delving into theory: color schemes based on HSV/HLC models generally just look good.



R can translate colors to rgb (this is handy for matching colors in other programs)

```
cat(paste("rgb", "#FF0000", ", blue"))
```

**R Color Palettes**

This is for all of you who don't know anything about color theory, and don't care but want some nice colors on your map or figure....NOW!

**TIP:** When it comes to selecting a color palette, DO NOT try to handpick individual colors! You will waste a lot of time and the result will probably not be all that great. R has some good packages for color palettes. Here are some of the options

Packages: `grDevices` and `colorspace`

`grDevices` comes with the base installation and `colorRamps` must be installed. Each palette's function has an argument for the number of colors and transparency (`alpha`):

grDevices palettes
<code>coltint</code>
<code>topo.colors</code>
<code>terrain.colors</code>
<code>heat.colors</code>
<code>rainbow</code>
<code>see p. 4 for options</code>

`heat.colors(4, alpha=1)`  
`> #FF0000FF "#FF0000FF" "#FFFF00FF" "#000000FF"`

For the rainbow palette you can also select start/end colors (`r = 0, end = 1, green = 1/6, green = 2/6, open = 3/6, blue = 4/6 and magenta = 5/6) and saturation (s) and value (v):`

```
rainbow(s = 1, v = 1, start = 0, end = max(s), n = 1), alpha = 1
```

**Package: RColorBrewer**

This function has an argument for the number of colors and the color palette (see P. 4 for options).

```
brewerpalettes[4, "Set3"]
```

```
> "#80002C" "#4DB6AC" "#BABA0A" "#800072"
```

To view colorbrewer palettes in R: `library(brewer.all)`

There is also a very nice interactive viewer: <http://colorbrewer.org/>

**## My Recommendation ##**

**Package: colorspace**

These color palettes are based on HCL and HSV color models. The results can be very aesthetically pleasing. There are some default palettes:

colorspace default palettes
<code>diverge_hsv</code>
<code>diverge_hcl</code>
<code>terrain_hcl</code>
<code>sequential_hcl</code>
<code>rainbow_hcl</code>

`"#E495A5" "#BAB005" "#99BEB1" "#CACAE2"`

However, all palettes are fully customizable.

`diverge_hcl(7, h = c(246, 40), c = 96, l = c(65, 90))`

Choosing the values would be daunting. But there are some pre-defined palettes in the `colorspace` documentation. There is also an interactive tool that can be used to obtain a customized palette. To start the tool:

```
pal <- choose_palette()
```

Page 1, Melanie Frazier

Week 3/3

STAT1003 Lecture 3

**colorRamps and grDevices**

colorRamps
<code>topo.colors</code>
<code>heat.colors</code>
<code>rainbow</code>
<code>greenRed</code>
<code>blueYellow</code>
<code>greenRedYellow</code>
<code>heatRamp</code>
<code>heatRamp2</code>
<code>heatRamp3</code>
<code>heatRamp4</code>
<code>heatRamp5</code>
<code>heatRamp6</code>
<code>heatRamp7</code>
<code>heatRamp8</code>
<code>heatRamp9</code>
<code>heatRamp10</code>
<code>heatRamp11</code>
<code>heatRamp12</code>
<code>heatRamp13</code>
<code>heatRamp14</code>
<code>heatRamp15</code>
<code>heatRamp16</code>
<code>heatRamp17</code>
<code>heatRamp18</code>
<code>heatRamp19</code>
<code>heatRamp20</code>
<code>heatRamp21</code>
<code>heatRamp22</code>
<code>heatRamp23</code>
<code>heatRamp24</code>
<code>heatRamp25</code>
<code>heatRamp26</code>
<code>heatRamp27</code>
<code>heatRamp28</code>
<code>heatRamp29</code>
<code>heatRamp30</code>
<code>heatRamp31</code>
<code>heatRamp32</code>
<code>heatRamp33</code>
<code>heatRamp34</code>
<code>heatRamp35</code>
<code>heatRamp36</code>
<code>heatRamp37</code>
<code>heatRamp38</code>
<code>heatRamp39</code>
<code>heatRamp40</code>
<code>heatRamp41</code>
<code>heatRamp42</code>
<code>heatRamp43</code>
<code>heatRamp44</code>
<code>heatRamp45</code>
<code>heatRamp46</code>
<code>heatRamp47</code>
<code>heatRamp48</code>
<code>heatRamp49</code>
<code>heatRamp50</code>
<code>heatRamp51</code>
<code>heatRamp52</code>
<code>heatRamp53</code>
<code>heatRamp54</code>
<code>heatRamp55</code>
<code>heatRamp56</code>
<code>heatRamp57</code>
<code>heatRamp58</code>
<code>heatRamp59</code>
<code>heatRamp60</code>
<code>heatRamp61</code>
<code>heatRamp62</code>
<code>heatRamp63</code>
<code>heatRamp64</code>
<code>heatRamp65</code>
<code>heatRamp66</code>
<code>heatRamp67</code>
<code>heatRamp68</code>
<code>heatRamp69</code>
<code>heatRamp70</code>
<code>heatRamp71</code>
<code>heatRamp72</code>
<code>heatRamp73</code>
<code>heatRamp74</code>
<code>heatRamp75</code>
<code>heatRamp76</code>
<code>heatRamp77</code>
<code>heatRamp78</code>
<code>heatRamp79</code>
<code>heatRamp80</code>
<code>heatRamp81</code>
<code>heatRamp82</code>
<code>heatRamp83</code>
<code>heatRamp84</code>
<code>heatRamp85</code>
<code>heatRamp86</code>
<code>heatRamp87</code>
<code>heatRamp88</code>
<code>heatRamp89</code>
<code>heatRamp90</code>
<code>heatRamp91</code>
<code>heatRamp92</code>
<code>heatRamp93</code>
<code>heatRamp94</code>
<code>heatRamp95</code>
<code>heatRamp96</code>
<code>heatRamp97</code>
<code>heatRamp98</code>
<code>heatRamp99</code>
<code>heatRamp100</code>

**colorspace defaults**

colorspace
<code>diverge_hsv</code>
<code>diverge_hcl</code>
<code>terrain_hcl</code>
<code>heat_hcl</code>
<code>sequential_hcl</code>
<code>rainbow_hcl</code>

**colorspace useful palette examples**

Example
<code>terrain_hcl(12, c = c(0, 55, 1 + c(45, 95, power = c(0.13, 1.5)))</code>
<code>heat_hcl(12, c = c(0, 30, 1 + c(30, 90, power = c(1.05, 1.5)))</code>
<code>heat_hcl(12, h = c(0, 100, 1 + c(75, 40, power = 1))</code>
<code>diverge_hcl(12, c = 100, 1 + c(50, 90, power = 1))</code>
<code>diverge_hcl(12, h = c(255, 330, 1 + c(40, 90)))</code>
<code>diverge_hcl(12, h = c(128, 330, t = 98, 1 - c(5, 90, 90)))</code>
<code>diverge_hcl(12, h = c(180, 330, t = 70, 1 - c(90, 90, 90)))</code>
<code>diverge_hcl(12, h = c(130, 45, c = 54, 1 + c(75, 95)))</code>
<code>diverge_hcl(12, h = c(180, 70, 1 + c(70, 90, 90)))</code>
<code>diverge_hcl(12, h = c(248, 40, c = 66))</code>

To begin interactive color selector: `pal <- choose_palette()`

**Useful Resources:**

A larger color chart of R named colors: <http://research.stowers-institute.org/efg/R/Color/Chart/ColorChart.pdf>

Nice overview of color in R: <http://research.stowers-institute.org/efg/Report/UsingColorInR.pdf>

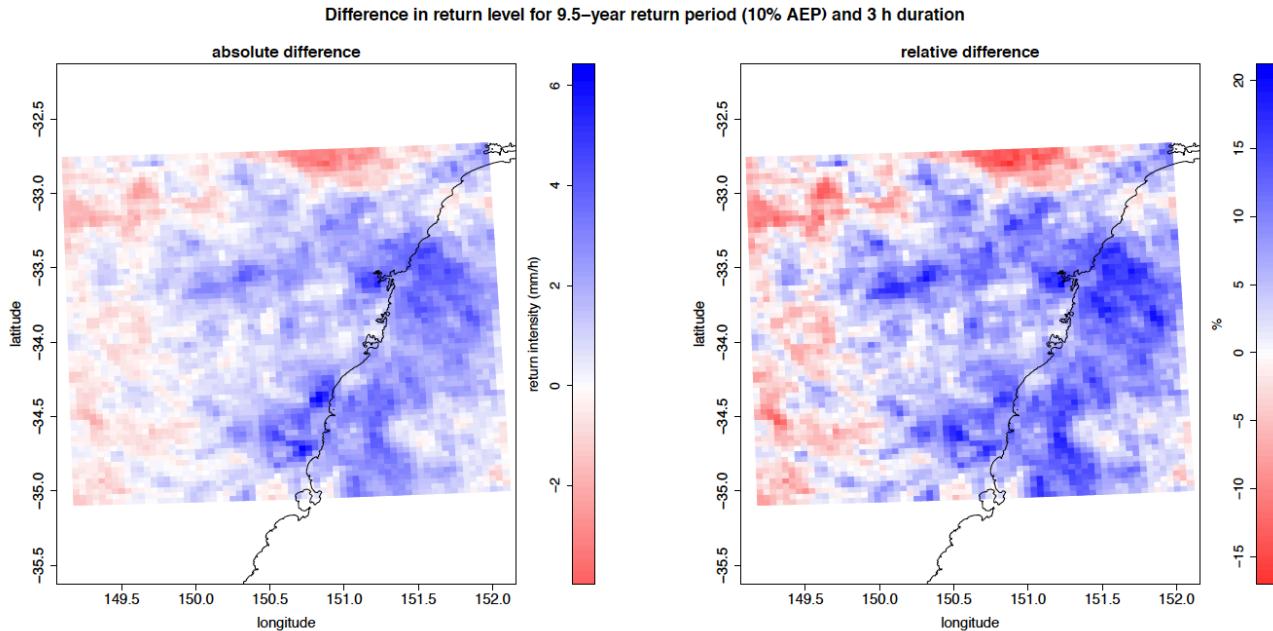
<http://students.washington.edu/mcklars0/documents/colorsVer2.pdf>

A color theory reference: Zelies, A. K. Hornik, P. Murrell. 2009. Escaping RColorland: selecting colors for statistical graphics. Computational and Statistics & Data Analysis 53:3259-3270

Page 4, Melanie Frazier

Curtin University

# Colour palette to highlight high and low values



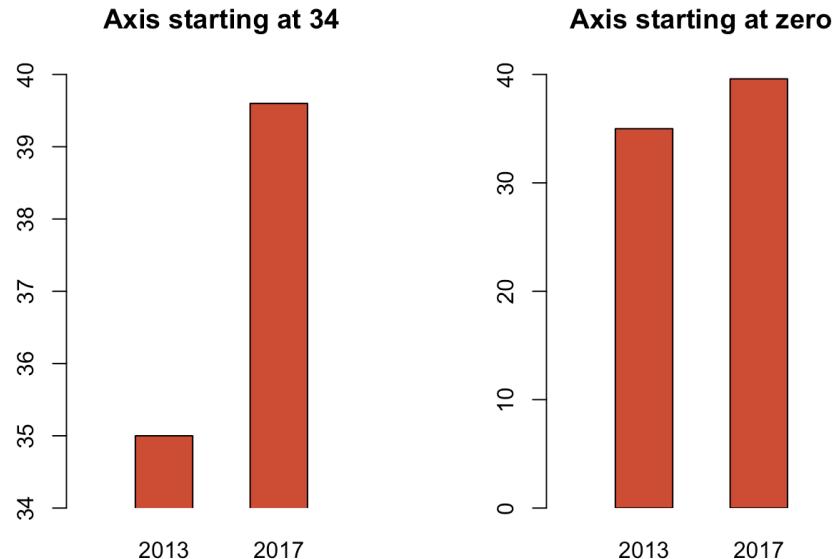
# One dataset, visualized 25 ways

- “Communication is in the message received”
  - It doesn’t matter what you intended, what counts is the viewer’s perception
- *How* we construct a plot can determine whether the viewer ‘gets’ the message we intended
- See <http://flowingdata.com/2017/01/24/one-dataset-visualized-25-ways/>

# Perception – Length

- We are pretty good at judging length visually, by comparison to area, volumes, or angles
- Practical consequence
  - Bar charts are to be preferred to pie charts
- Practical considerations
  - Visually, the longer a bar, the greater the absolute value it represents
  - Can be vertical or horizontal
  - To judge magnitude, bars must represent the entire length, not just a portion

# Spot the difference



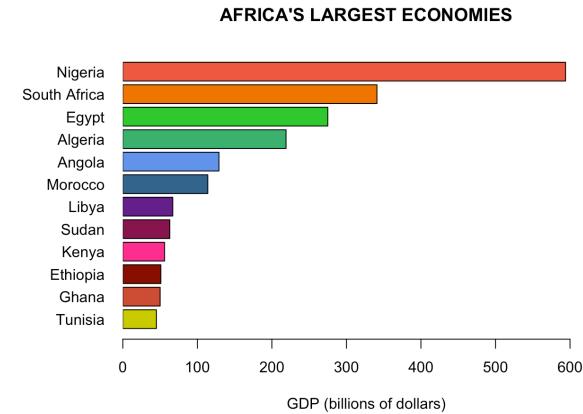
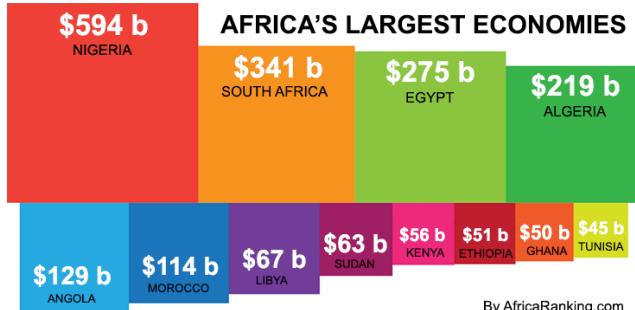
Week 3/3

STAT1003 Lecture 3

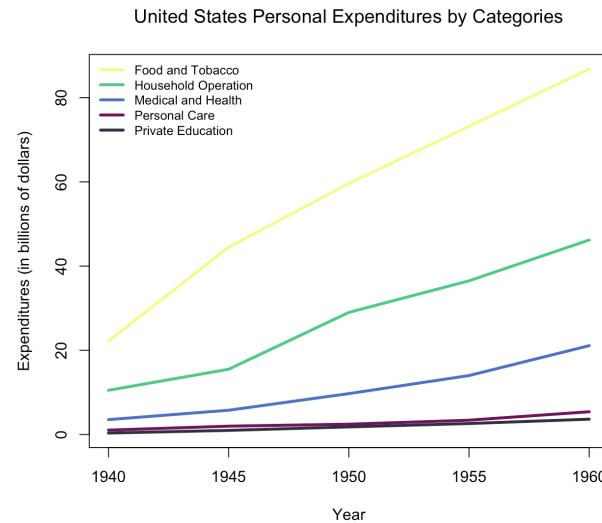
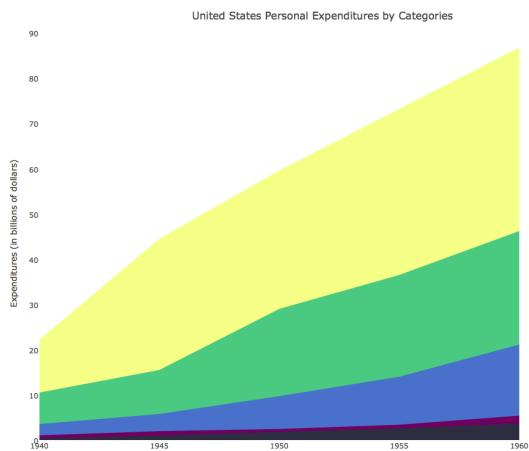


Curtin University

# Areas or lengths?



# Stacked barplots, histograms, and line plots



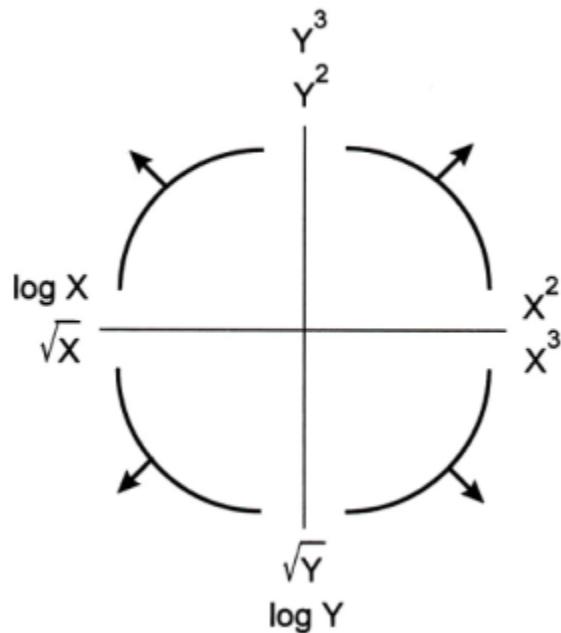
Week 3/3

STAT1003 Lecture 3

# Transformations

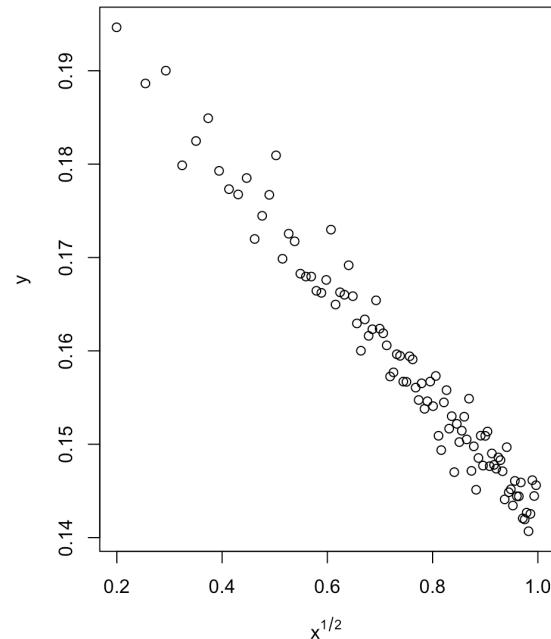
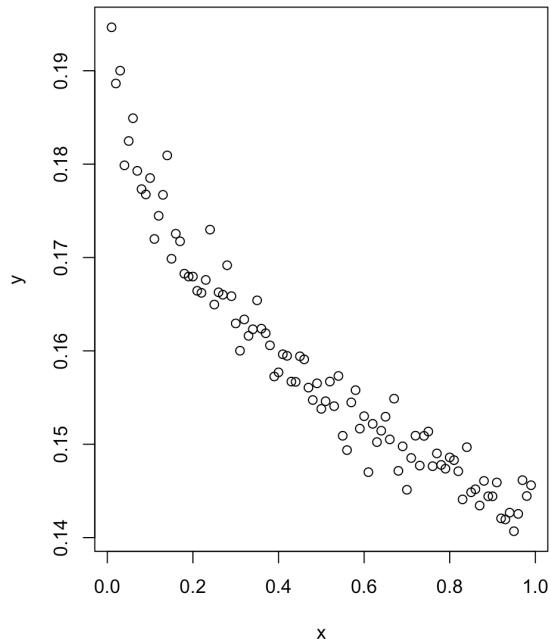
- Sometimes, curvilinear bivariate relationships can be straightened out by transforming one or more variables
  - Easier to interpret than nonlinear relationship
  - Prior to regression
- Transformations can be useful when
  - the data covers several orders of magnitude
  - data is ‘squashed in’

# Tukey-Mosteller's Bulging Rule



- For example, if you have data in the first quadrant, try transformations along the positive  $Y$ -axis or  $X$ -axis or both
- [http://www.dmstat1.com  
/res/EDATukeysBulgingRule.html](http://www.dmstat1.com/res/EDATukeysBulgingRule.html)

# Tukey-Mosteller's Bulging Rule

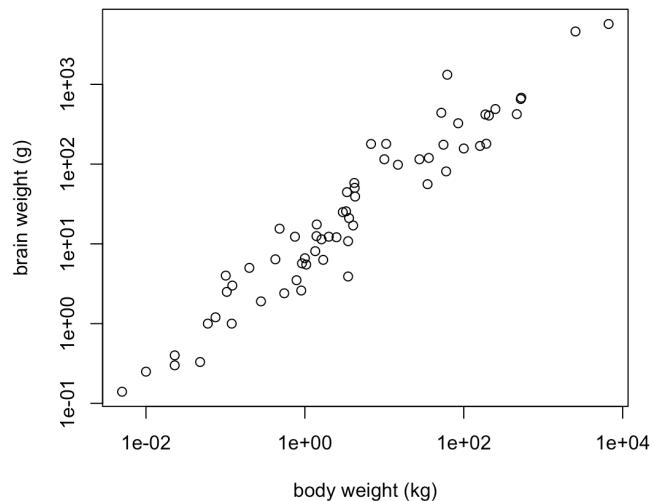
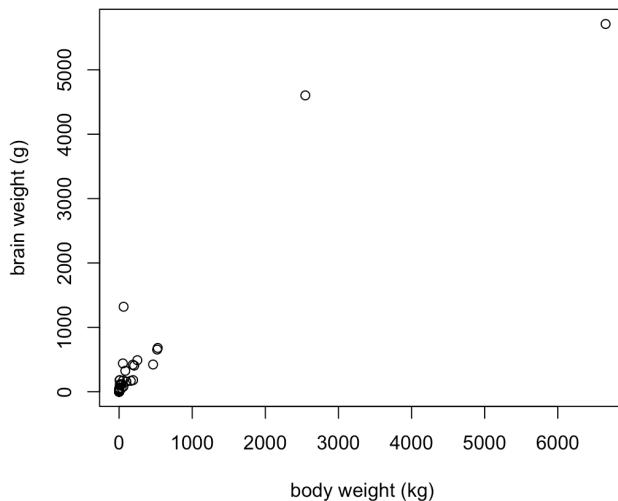


Week 3/3

STAT1003 Lecture 3

# Transformations to ‘zoom in’

Brain and body weights of mammals



# Context

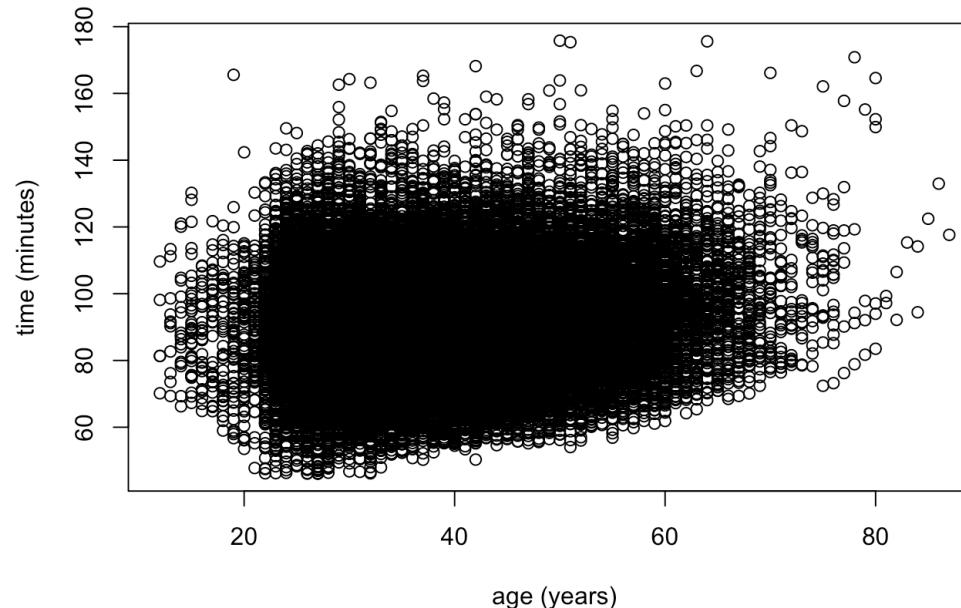
- Context is information that leads to a better understanding of the who, what, where, and why of your graphic, e.g.,
  - Axes labels
  - Units of measurements
  - Label points of unusual interest, e.g., humans in the previous plot
  - Depending on broader context in which the graphic is appearing, a caption describing the data and important features may be necessary

# Plotting lots of data

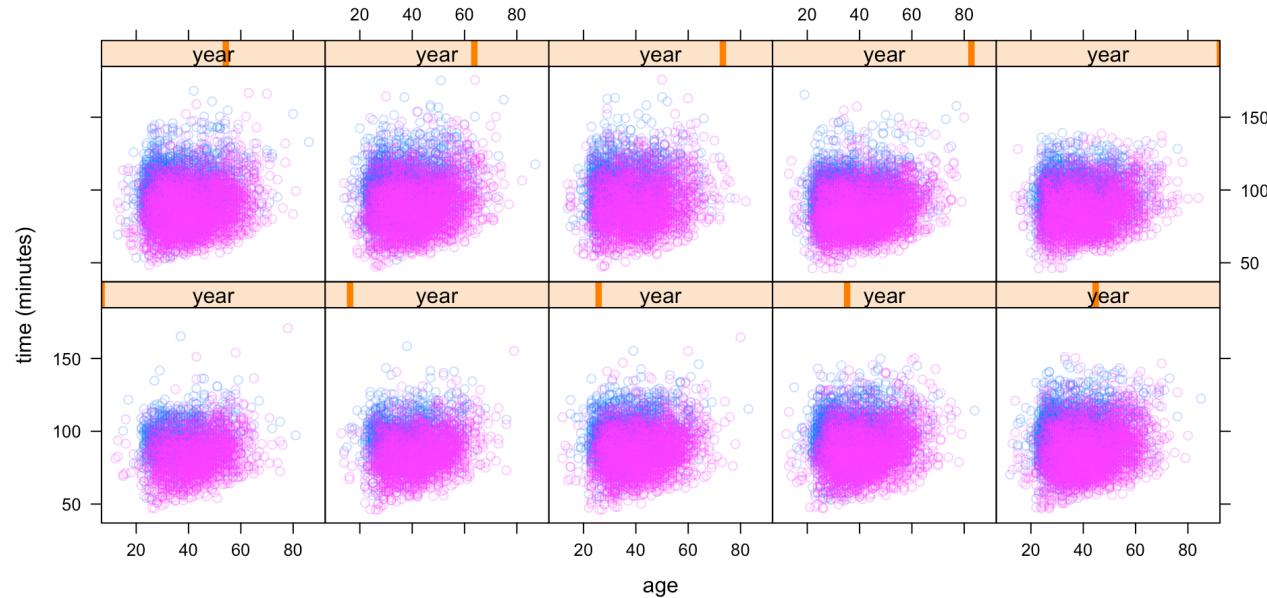
- Split the data into natural subgroups (conditioning) before plotting
- Use transparent colours
- Consider adding a scatterplot smoother to help guide the viewer's eye

# Cherry Blossom Run (Washington, DC)

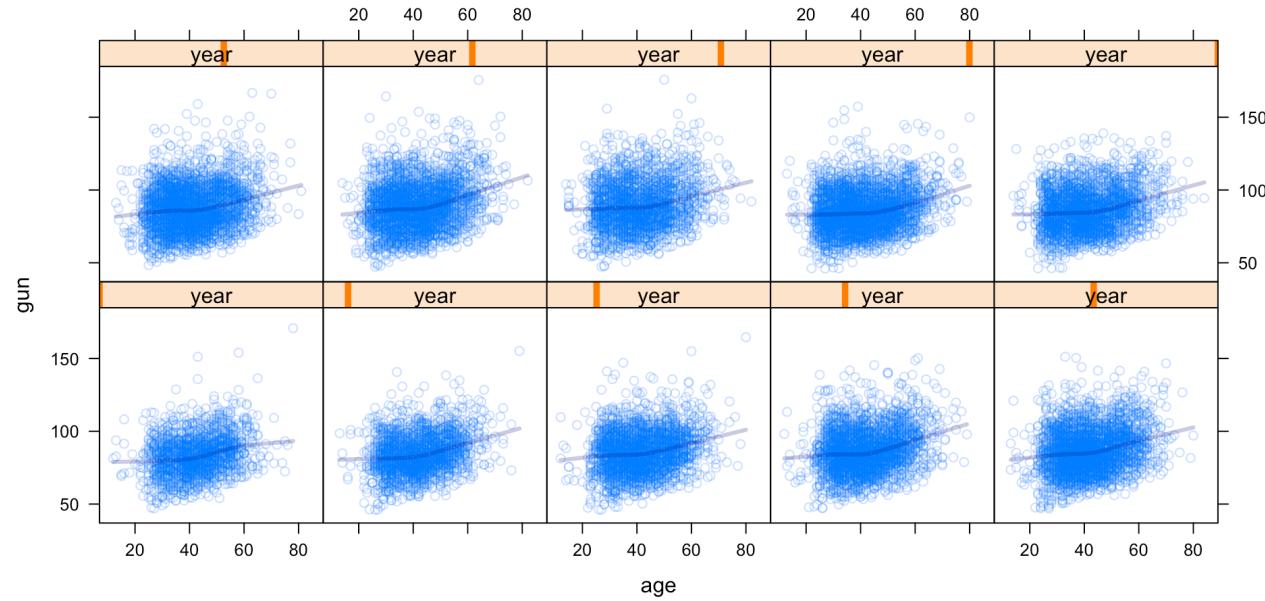
1999 - 2008: Over 40000 runners!



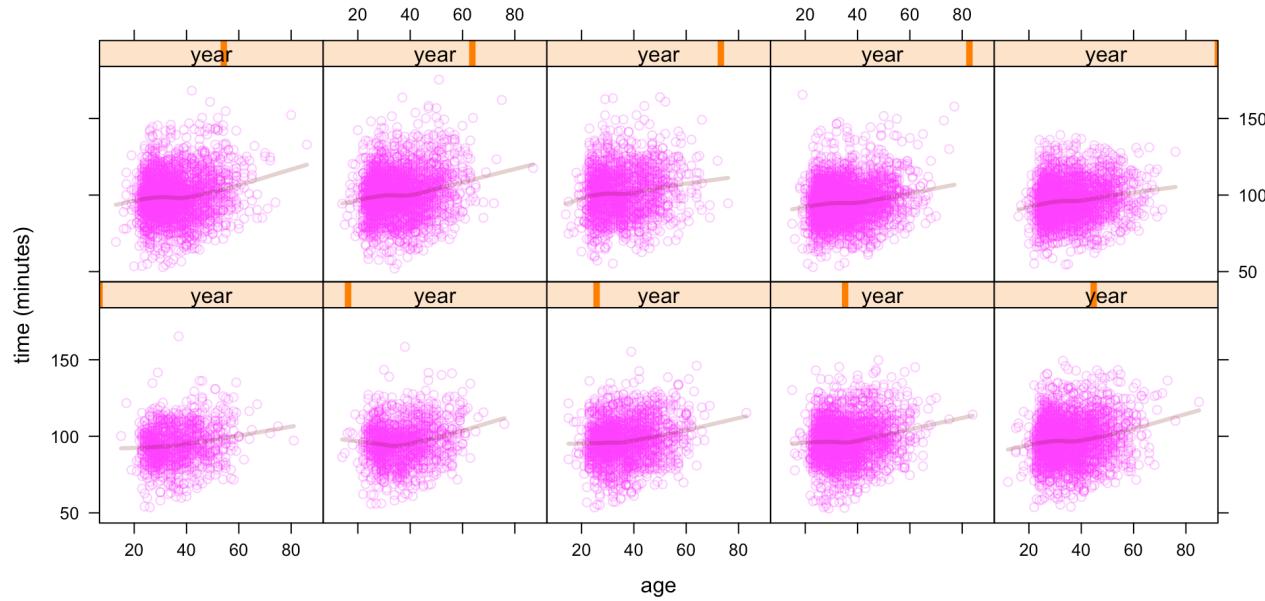
# Cherry Blossom Run: conditioning & transparency



# Cherry Blossom Run: adding a scatterplot smoother (males)

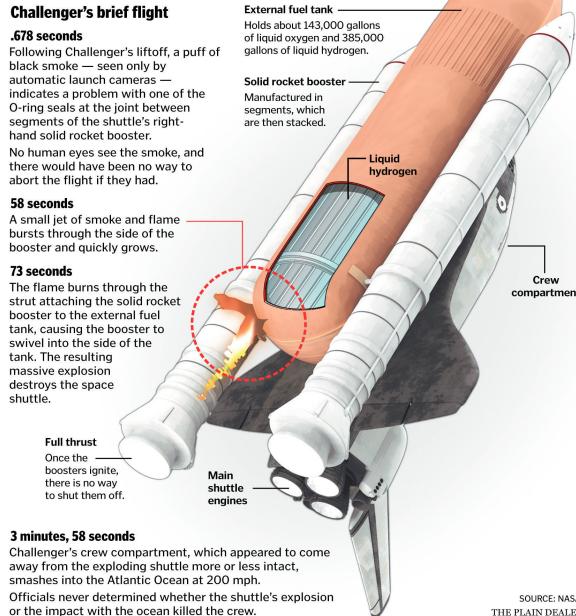


# Cherry Blossom Run: adding a scatterplot smoother (females)



# The importance of good visualization: the *Challenger* disaster (1986)

## A major malfunction



Week 3/3

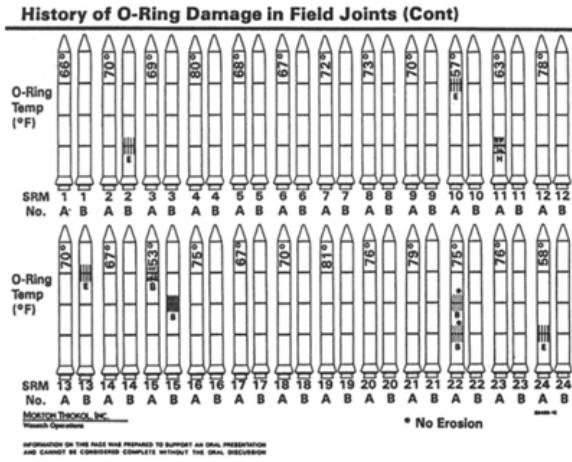
STAT1003 Lecture 3

# *Challenger* (Baumer *et al.* 2017)

- Engineers at Morton Thiokol, who supplied solid rocket motors to NASA, recommended that the launch be delayed because cold weather would jeopardize the stability of O-rings
- Their recommendations were overruled due to lack of persuasive evidence, and 73 seconds after launch, *Challenger* exploded
- The evidence was in the form of hand-written tables, but none was *graphical*, and they failed to convince NASA to delay the launch

# Challenger

MOTOR	HISTORY OF O-RING TEMPERATURES (DEGREES - F)			
	MST	AMB	O-RING	WIND
DM-4	68	36	47	10 MPH
DM-2	76	45	52	10 MPH
QM-3	72.5	40	48	10 MPH
QM-4	76	48	51	10 MPH
SRM-15	52	64	53	10 MPH
SRM-22	77	78	75	10 MPH
SRM-25	55	26	29	10 MPH
			27	25 MPH



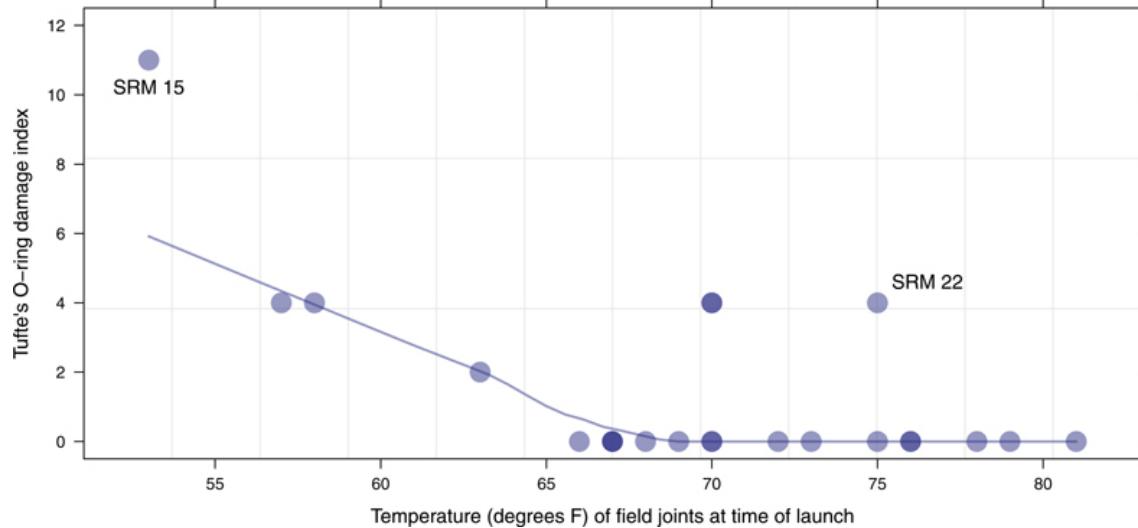
- (a) One of the original 13 charts presented by Morton Thiokol engineers to NASA on the conference call the night before the Challenger launch. This is one of the more data-intensive hearings after the Challenger explosion. This is charts.
- (b) Evidence presented during the congressional hearing. This is a classic example of “chartjunk.”

# *Challenger*

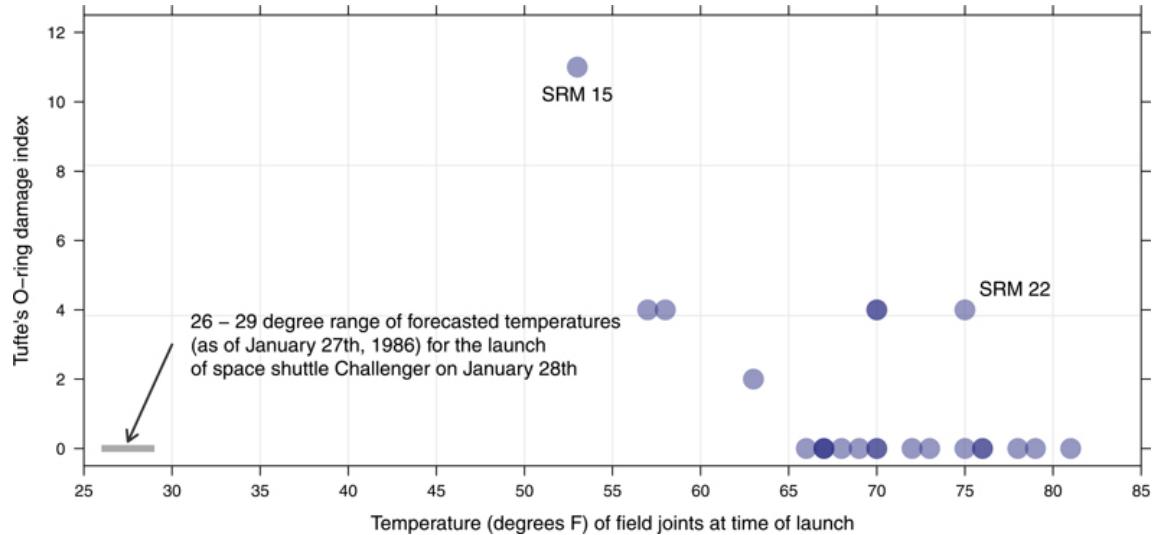
## Critiques by Tufte (1997)

- The tables (charts) that were discussed at meetings made it difficult to see bivariate (temperature/failure) relationships
- Anecdotal evidence
  - With small sample size, anecdotal evidence can be particularly challenging to refute
  - Engineers argued that SRM-15 had the most damage on the coldest previous launch date, NASA officials countered that SRM-22 had almost as much damage on one of the warmer launch dates
- *Plotting the data tells a different story!*

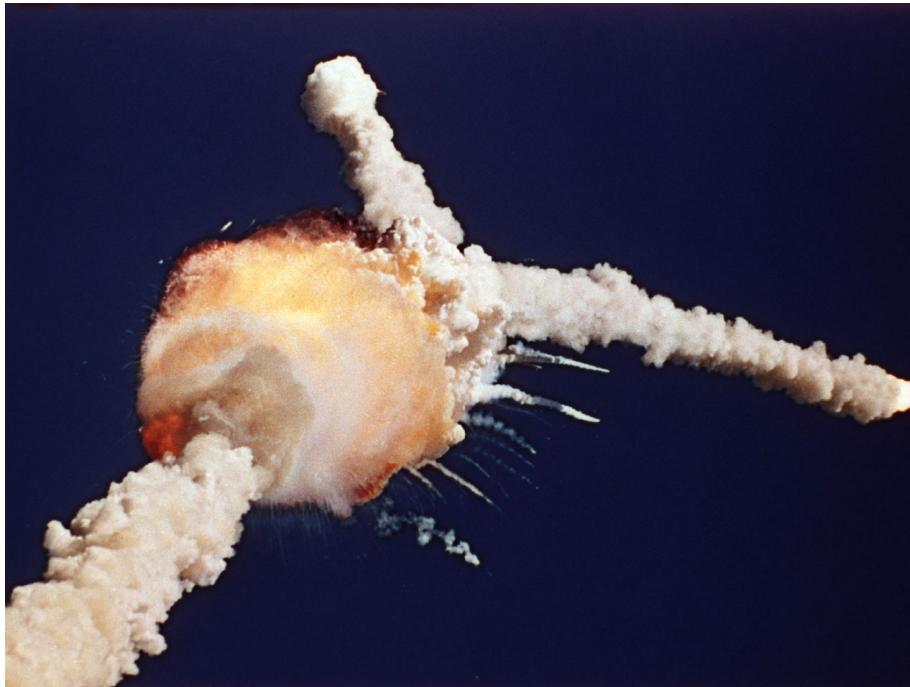
# Tufte (1997)



# Tufte (1997)



# *Challenger (1986)*



Week 3/3

STAT1003 Lecture 3



# Guidelines (Nolan, 2017)

Reveal the data

- Choose scale appropriately
- Avoid having other graph elements interfere with the data
- Use visually prominent symbols
- Eliminate superfluous material – ‘chart junk’ – as in the *Challenger* graphic
- Avoid plotting points on top of one another - use jittering

# Guidelines (continued)

- Put juxtaposed plots on the same scale
- Make it easy to distinguish elements of superposed plots, e.g., colour, line type
- Avoid stacking
- Avoid areas, volumes
- Don't break visual metaphor, i.e., if using rectangles, then area should correspond to the value of the variable you're trying to represent, not the length of one of the sides

# Guidelines (continued)

- Describe what you want the reader to see in the caption
- Use informative labels and legends
- Use colour and plotting symbols to add more context
- Plot the same thing in many ways/scales, but not on the same plot!

# Taxonomy of plotting methods (Nolan, 2017)

Type	Plot
Numeric	Few observations Histogram, Density curve Box plot, Violin plot Normal quantile plot Few Observations - Rug plot, Dot plot Caution if discrete: density curves and box plots may be misleading
Categorical	Dot chart Bar chart Pie chart (avoid!) Caution if ordinal – order of bars, dots, etc. should reflect category order

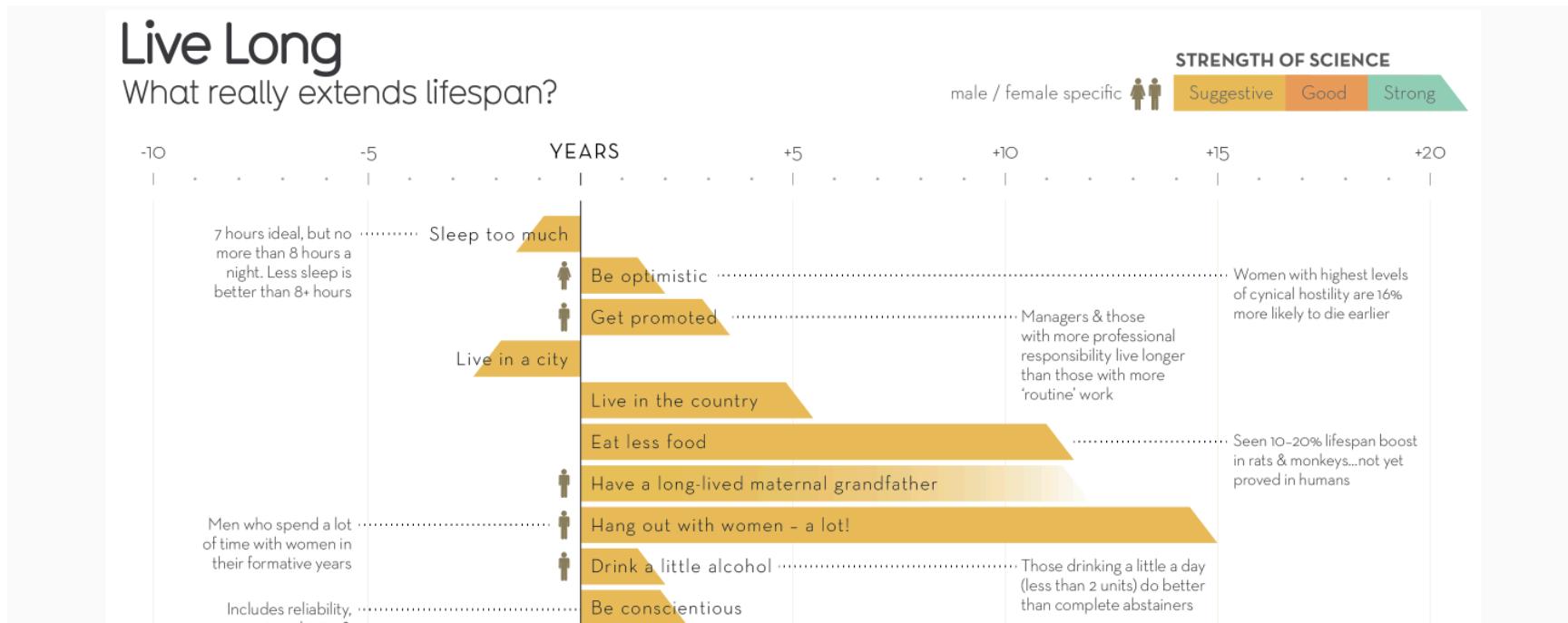
# Taxonomy of plotting methods (continued)

	Numeric	Categorical
Numeric	Scatter plot Smooth scatter Smooth lines and curves	Multiple histograms, density curves, Avoid jiggling baselines
Categorical		Side-by-side bar plot Overlaid lines plot Side-by-side dot chart Mosaic plot Avoid stacking

# More misleading visualizations

- Internet Explorer and murders in US
- GHG emissions
- US presidents and the unemployment rate
- Age structure of college enrollments
- Yet more misleading visualizations...

# Misleading visualization I



Week 3/3

STAT1003 Lecture 3



Curtin University

# Next week

- Exploratory data analysis

