

# DATA VISUALIZATION AND R

## THEORY AND IMPLEMENTATION

### SPRING 2016

Ryan Womack([rwmack@rutgers.edu](mailto:rwmack@rutgers.edu))  
Data Librarian, Rutgers University

March 31, 2016



This work is licensed under a [Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License](#).

# INTRODUCTION

What this workshop IS:

- Focuses on standard techniques of data visualization, the day-to-day power tools for understanding data
- Reviews various graphical techniques, from early to recent, from simple to advanced
- Discusses principles of good data presentation, and show the R implementation of many functions

# INTRODUCTION

What this workshop is NOT:

- It is not about “infographics”, the beautiful, heavily customized products of expert graphic designers
- Not about the cognitive science aspects of data perception [wish I knew more about this!]
- It is not a complete introduction to R, even though R is used
- It is not an introduction to other software beyond R, or the use of R with scripting languages to produce interactive graphics on the web
- It is not necessarily a balanced survey of all data visualization. In particular, it is light on graph networks, clustering, and trees [not my expertise]
- Very little mapping, too. [Done to death?]

# INTRODUCTION

What you can hope to gain:

- Familiarity with the basic principles and history of data visualization, and recent developments
- Exposure to a wide-range of plotting techniques and R packages
- A sampling of interactive and big data methods
- Understanding of the power and potential of combining appropriate graphics techniques with data

# SETUP

- Workshop materials, including R scripts, supplemental images and data, are available for download from  
<https://github.com/ryandata/DataViz>
- The R script file contains working demonstrations of the concepts mentioned here.
- You will have to install any packages not already on your system.

# WHY DATA VISUALIZATION?

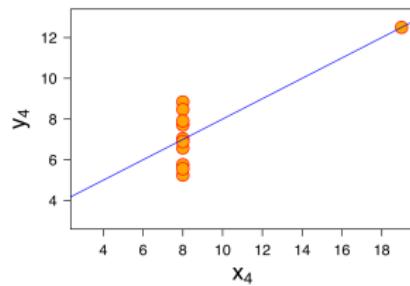
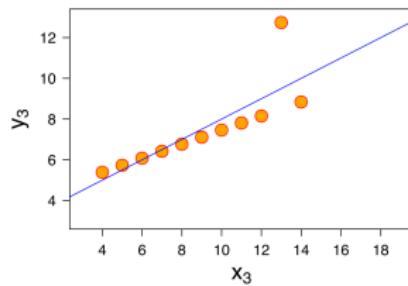
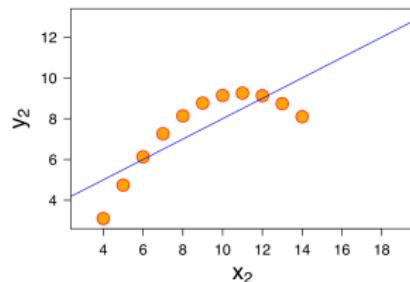
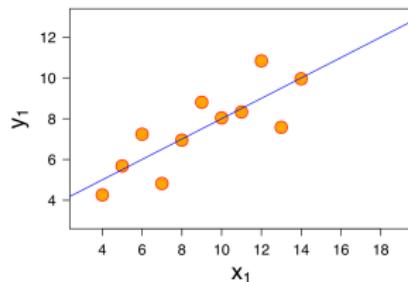
Data visualization can:

- provide clear understanding of patterns in data
- detect hidden structures in data
- condense information

# ANScombe'S QUARTET

For example, see Anscombe's quartet (image source:

[http://commons.wikimedia.org/wiki/File:Anscombe%27s\\_quartet\\_3.svg](http://commons.wikimedia.org/wiki/File:Anscombe%27s_quartet_3.svg)):



# LINKS TO DATAVIZ SITES

Some examples of good data visualization (and fancy infographics) can be found at:

- [Information Aesthetics](#)
- [Chart Porn](#)
- [Eagereyes](#)
- [DataVis.ca](#)
- [Visualizing.org](#)
- [VizWiz](#)
- [US Census Data Visualization Gallery](#)

# BAD GRAPHS

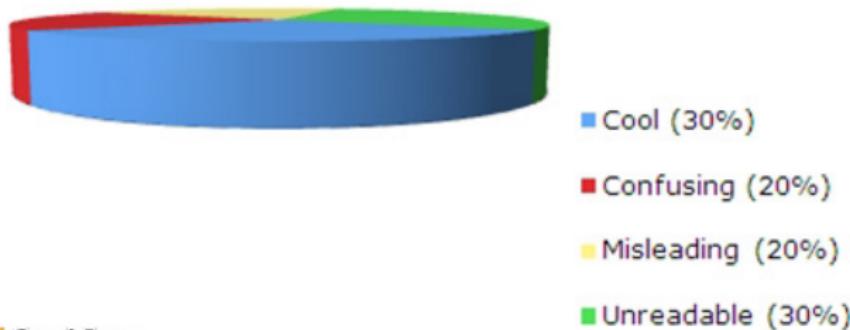
- Pie Charts are known to be problematic
- Clutter and other issues can ruin graphics

For more bad ideas, try:

- Junk Charts
- Ten Worst Graphs
- WTFviz

# PIE CHART EXAMPLES

## Perception of 3D pie charts



GraphJam

image source: <http://peltiertech.com/WordPress/3d-pie-charts/>

# PIE CHART EXAMPLES

Microsoft Word Features By Version Added

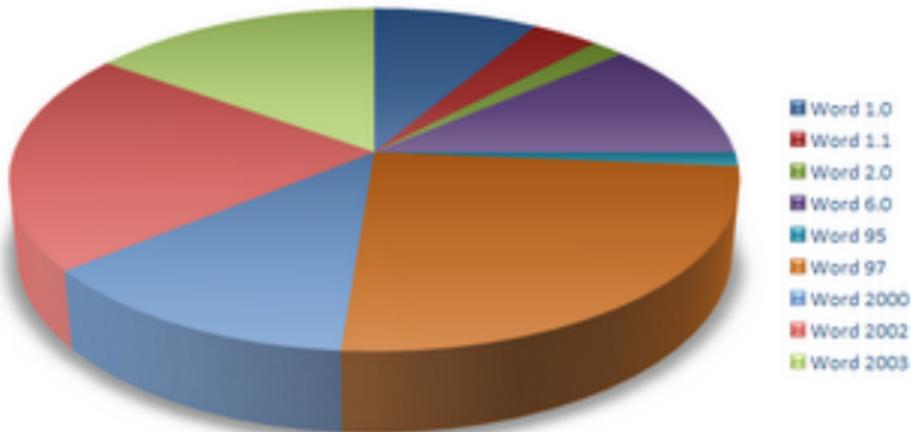


image source: <http://ndevisual.wordpress.com/tag/uses-of-pie-charts/>

# PIE CHART EXAMPLES

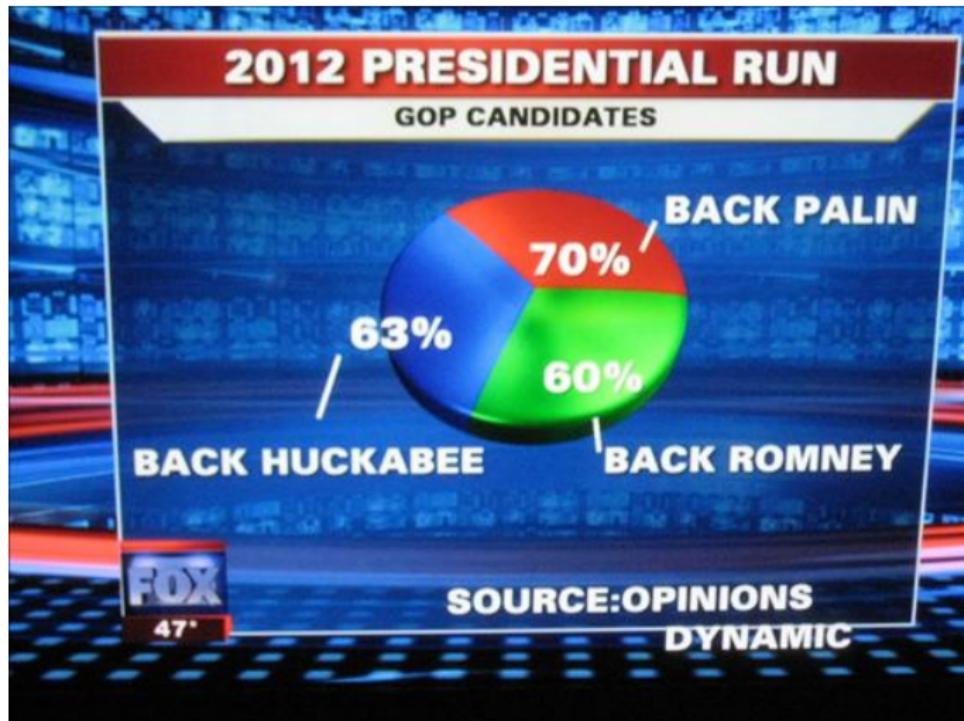


image source: <http://www.nbcchicago.com/news/local/FOX-News-Chart-Fails-Math-73711092.html>

# PIE CHART EXAMPLES

## Pie Charts Aren't Accurately Visually Interpreted

Question 1



Question 2



Question 3



- Product A ■ Product B
- Product C ■ Product D
- Product E

vovici.com

image source: [http://tips.vovici.com/content/111031\\_swb](http://tips.vovici.com/content/111031_swb)

# PIE CHART EXAMPLES

Bar Charts are Easier to Interpret than Pie Charts

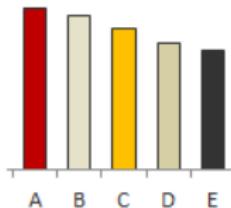
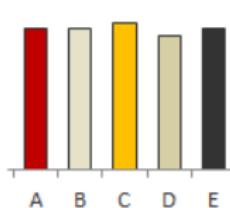
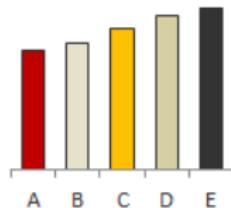
Question 1



Question 2



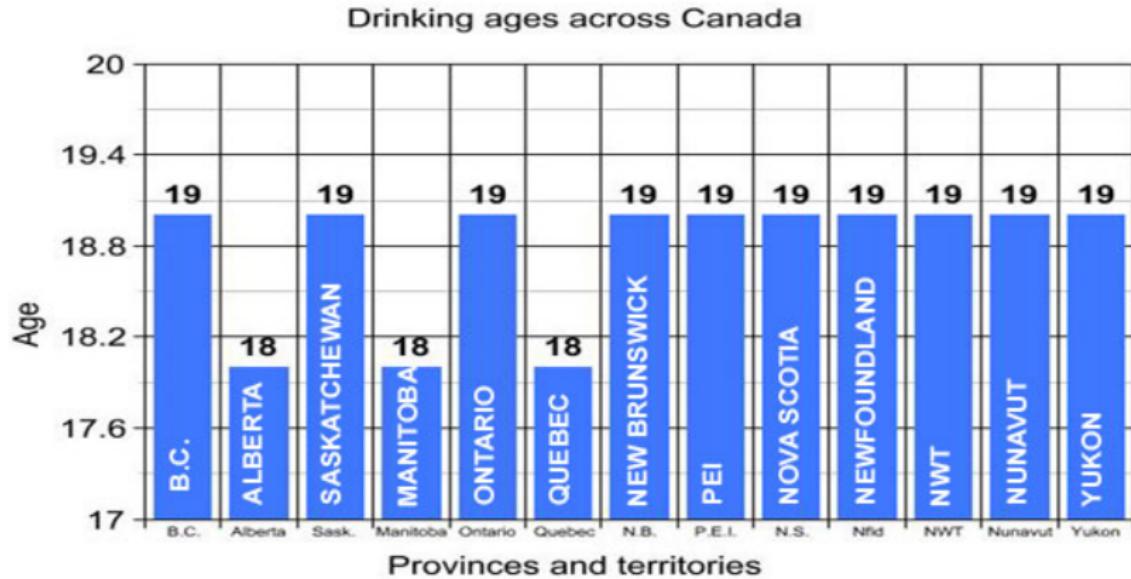
Question 3



vovici.com

image source: [http://tips.vovici.com/content/111031\\_swb](http://tips.vovici.com/content/111031_swb)

# CLUTTER EXAMPLE



Canadian Centre on Substance Abuse

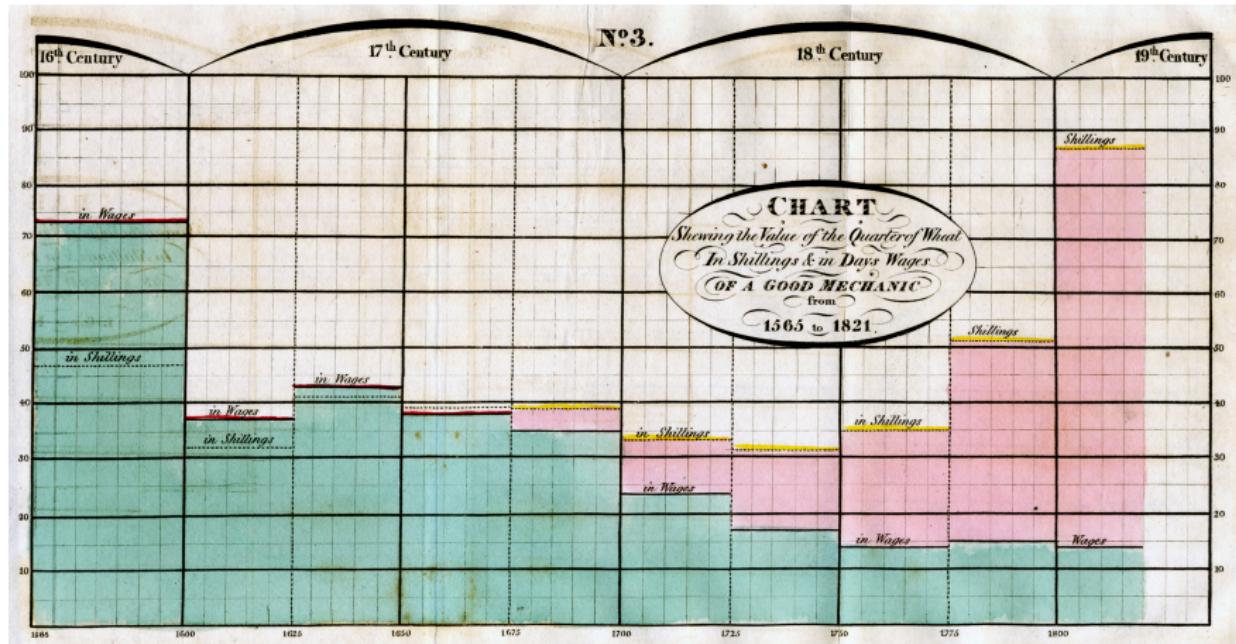
image source:

[http://junkcharts.typepad.com/junk\\_charts/2013/03/which-software-is-responsible-for-this.html](http://junkcharts.typepad.com/junk_charts/2013/03/which-software-is-responsible-for-this.html)

# PLAYFAIR

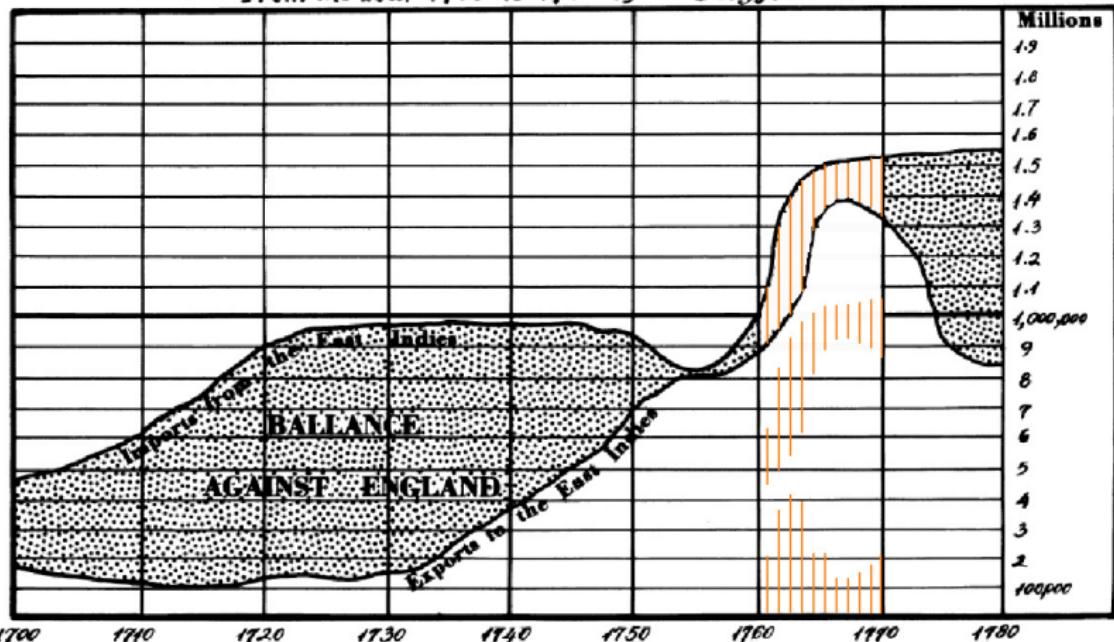
- Astronomical observations, charts, and maps led in graphical innovation prior to 1800.
- William Playfair is the pioneer of the line chart, bar chart, time series plots, and pie chart.
- Playfair, W. (1786). *Commercial and Political Atlas: Representing, by Copper-Plate Charts, the Progress of the Commerce, Revenues, Expenditure, and Debts of England, during the Whole of the Eighteenth Century,*
- Playfair, W. (1801). *Statistical Breviary.*
- Both republished in *The Commercial and Political Atlas and Statistical Breviary*, 2005, Cambridge University Press.

# PLAYFAIR EXAMPLES



# PLAYFAIR EXAMPLES

CHART of EXPORTS and IMPORTS to and from the EAST INDIES  
From the Year 1700 to 1780 by W. Playfair

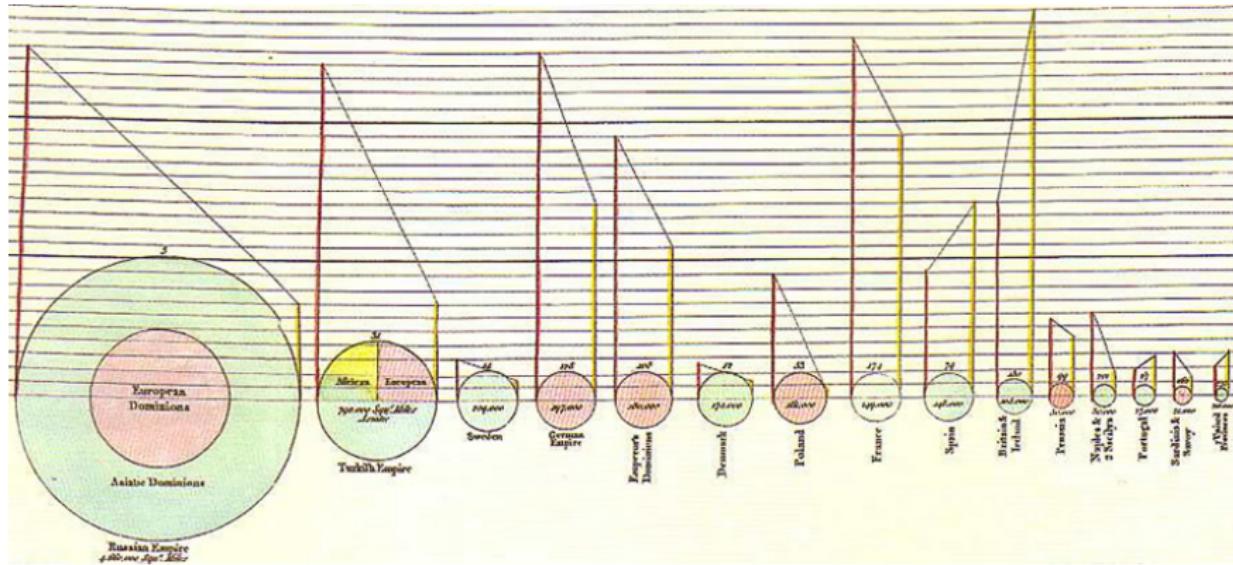


The Bottom Line is Divided into Years the Right-hand Line into HUNDRED THOUSAND POUNDS  
See page 34<sup>th</sup>

Published in the Act Divine 16<sup>th</sup> Aug. 1785



# PLAYFAIR EXAMPLES



# PIKETTY AND PLAYFAIR?

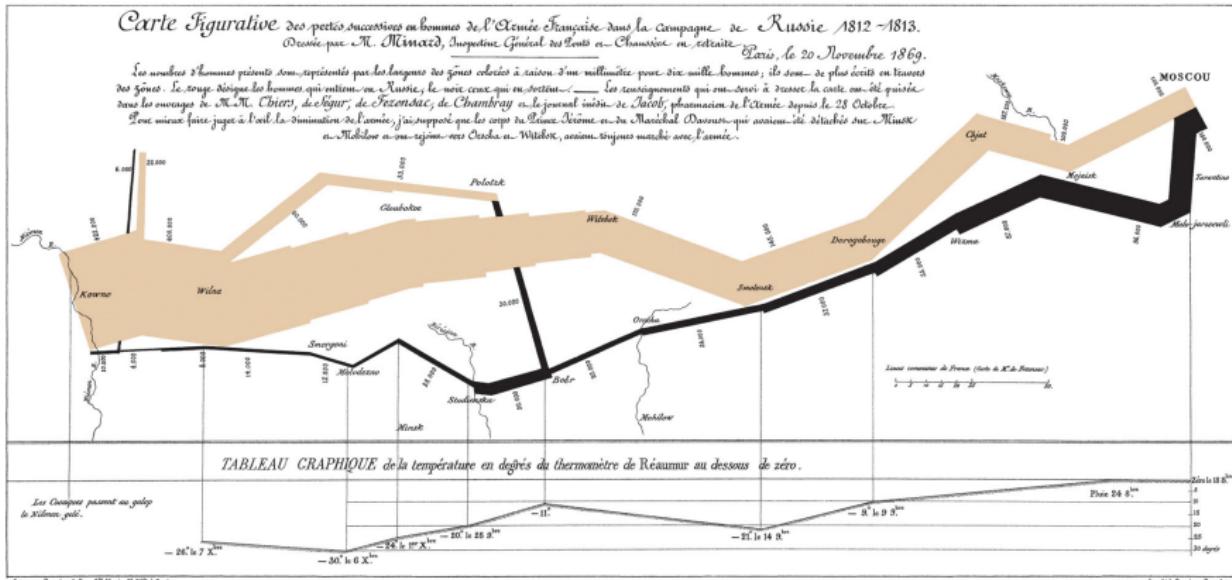
- Significance magazine suggests that Piketty's Capital garnered attention partly because of its careful use of figures in a style similar to Playfair's.
- But also that his charts could use some polishing!

# MINARD

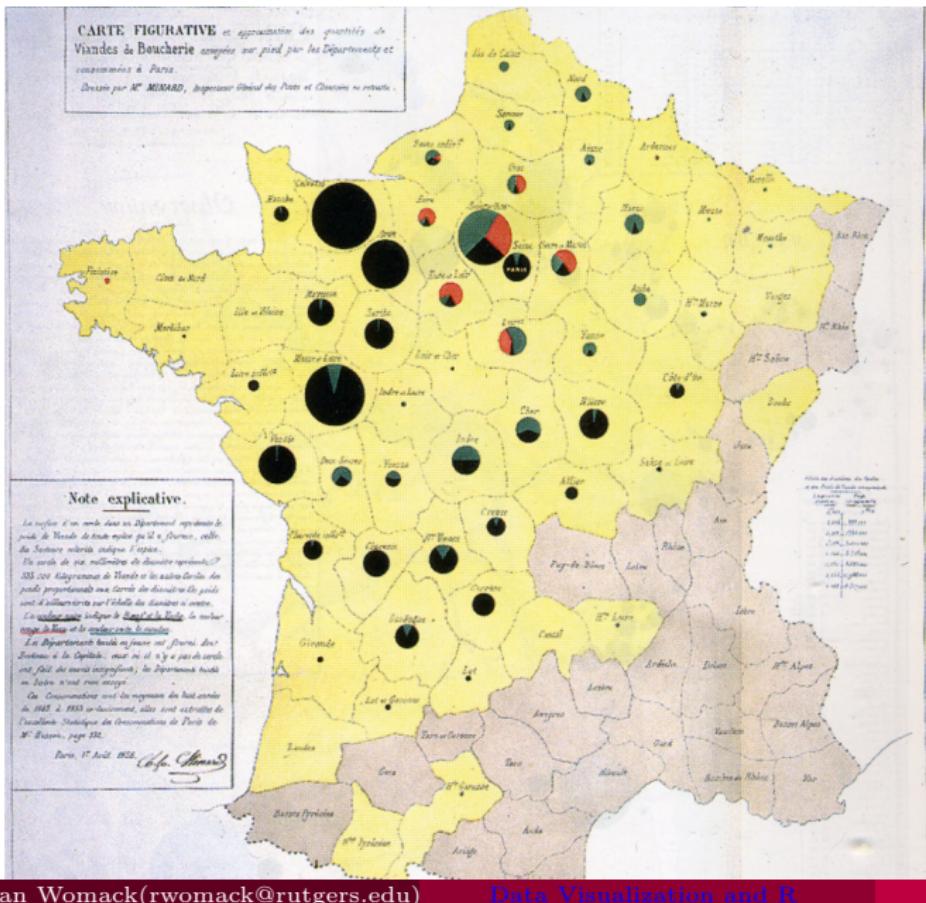
Charles Joseph Minard was the next influential data graphic creator after Playfair.

- Minard's [flow map of Napoleon's Russian campaign](#) is celebrated by Tufte and others as one of the greatest information graphics.
- It embodies an ideal of highly compressed informative elements, presented with style
- Six variables: size, location in 2 dimensions, the direction of the army, temperature, date [and group]
- However, this is a one-off design that crosses into Infographics, but it can be reproduced in [R](#) and other software.

# MINARD EXAMPLES



# MINARD EXAMPLES



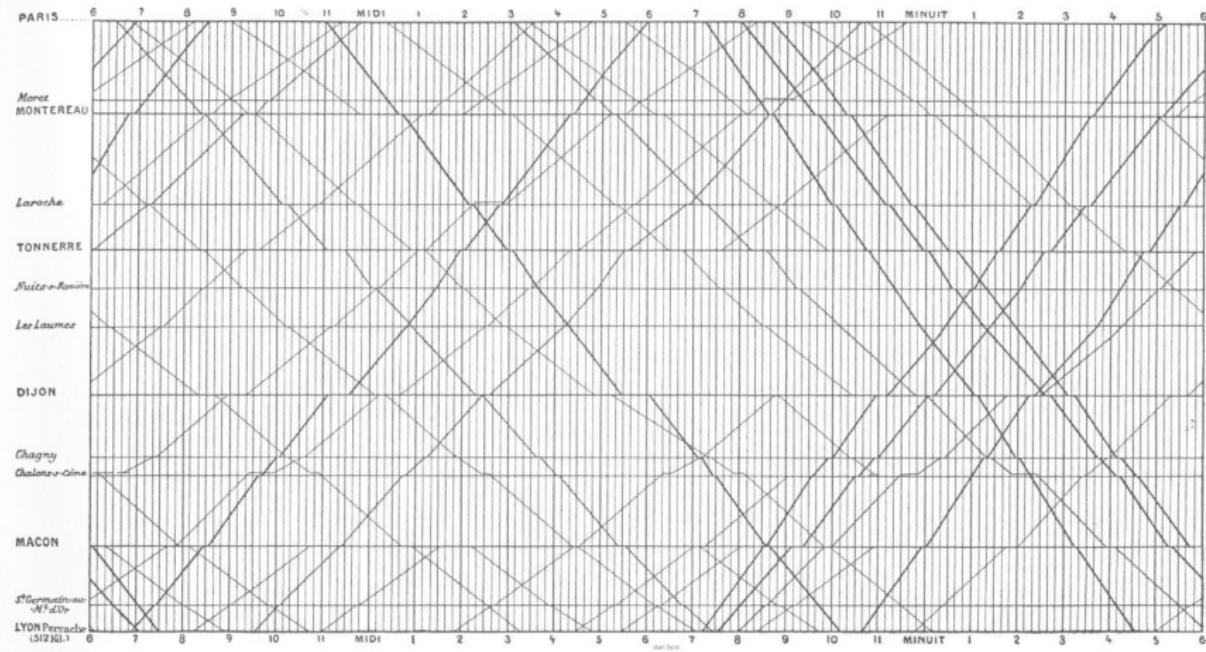
## FISHER AND TUKEY

- Statisticians such as Ronald Fisher and John Tukey continued to advance graphical methods for the analysis of data.
- Fisher emphasized plotting the data to understand relationships.
- Tukey's *Exploratory Data Analysis* emphasized the use of graphics to understand the data during analysis, rather than the final presentation to an outside audience.
- Tukey created the box and whiskers plot and the stem and leaf plot.

Edward R. Tufte's series of books, beginning with *The Visual Display of Quantitative Information*, have become the most widely known works on data visualization.

- There is considerable overlap between the various publications
- Tufte's ideal is highly compressed, elegant, and informative data, as expressed in dense printed graphics
- Tufte sometimes emphasizes beauty and design to the detriment of simplicity and clarity [e.g., train schedules]
- “Graphical elegance is often found in simplicity of design and complexity of data.”
- “Beautiful graphics do not traffic with the trivial.”

# TRAIN SCHEDULE FROM MAREY



# TUFTE'S PRINCIPLES

Tufte has developed and popularized numerous principles and terminology:

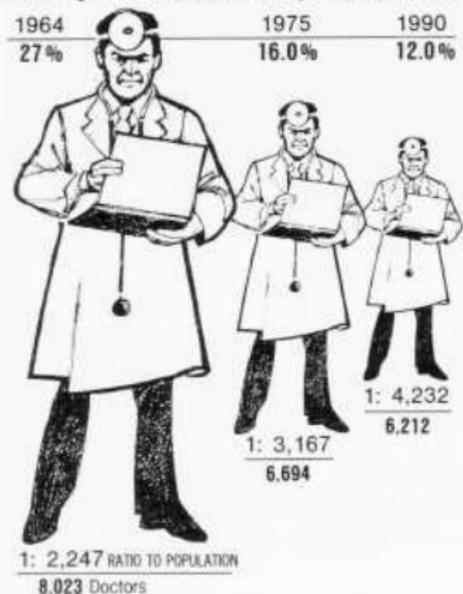
- **Graphics reveal data** - show the data without distorting it - “above all else show the data”
- **Small multiple** - understanding one slice makes understanding others easier
- **Lie factor** - effect shown/effect in reality
- **Graphical Integrity** - no lies, let data vary, not design
- **Data density** - maximize data/ink ratio
- **Sparklines** - seems they haven't caught on
- **chartjunk** - self-explanatory
- **Powerpoint** is responsible for most of the world's sorrows [*The Cognitive Style of Powerpoint*]

# LIE FACTOR

## THE SHRINKING FAMILY DOCTOR In California

Percentage of Doctors Devoted Solely to Family Practice

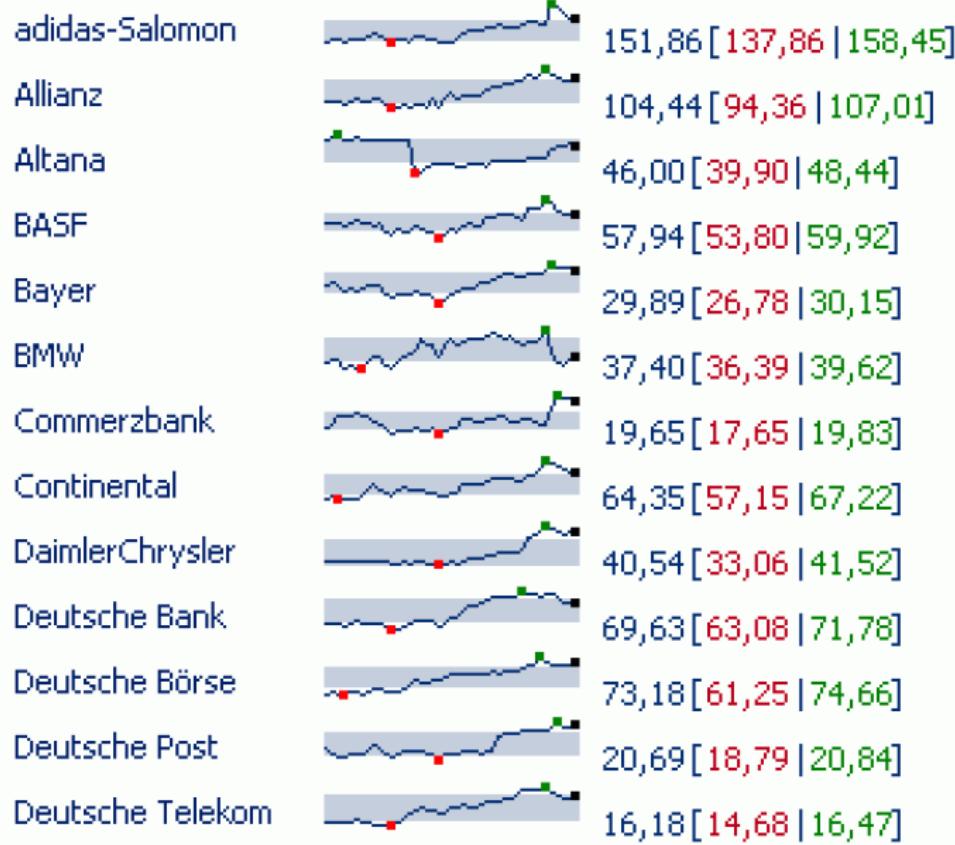
1964	1975	1990
27%	16.0%	12.0%



*Los Angeles Times*, August 5, 1979, p. J-

image source: <http://www.datavis.ca/gallery/lie-factor.php>

# TUFTE SPARKLINES



# TUFTÉ COUNTEREXAMPLE - MINARD REVISITED

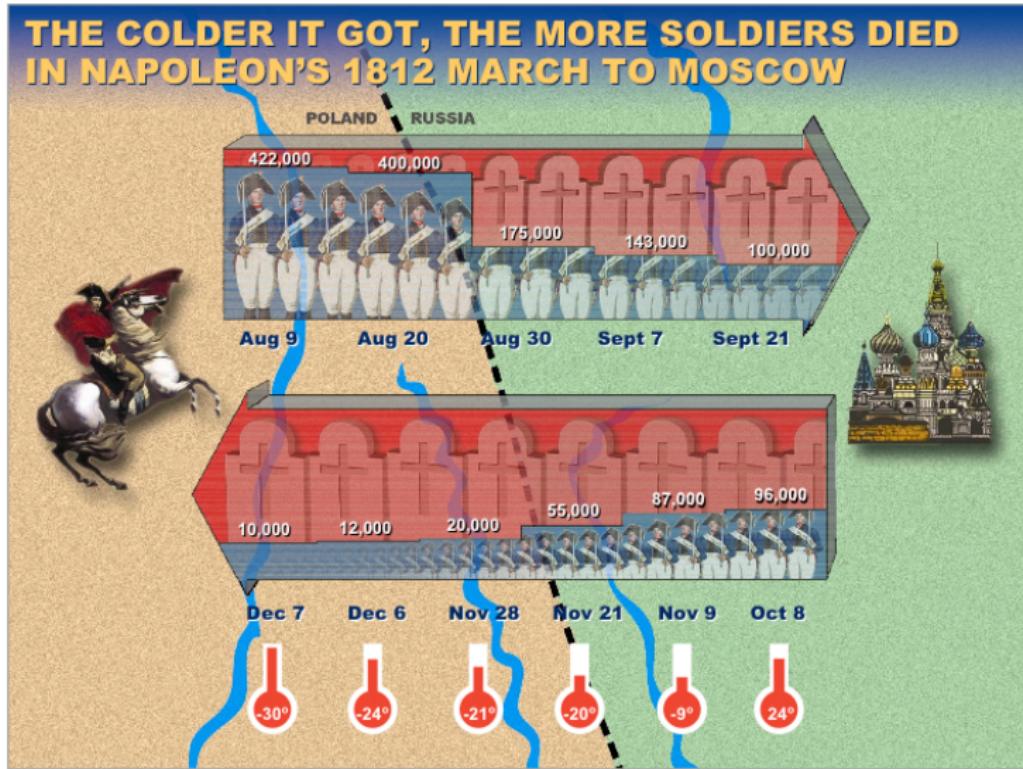


image source: Re-visions of Minard at <http://www.datavis.ca/gallery/re-minard.php>

# BACK TO PIES

Why is the pie chart bad?

- Low data density
- Failure to order numbers along a visual dimension
- Perception difficulty in judging area
  - Stacked bar charts also pose perceptual problems

Or in the terms of [Gary Klass](#):

- Data Ambiguity
- Data Distortion
- Data Distraction

# A GOOD PIE CHART?



image source:

<http://sciencegeekgirl.wordpress.com/2008/11/14/a-true-pie-chart/>

# CLEVELAND

- William Cleveland's *Elements of Graphing Data* and *Visualizing Data* pioneered systematic considerations of data legibility
- Cleveland is particularly known for promoting the *dot plot* as an alternative to bars and pies.
- The dot plot provides clarity and easy comparison of data.
- Cleveland also pioneered Trellis graphics
- Trellis graphics emphasizes comparison of multiple panels of data
- The `lattice` package implements Trellis graphics in R
- See [Cleveland.pdf](#) for a summary of Cleveland's recommendations

# CLEVELAND - TECHNIQUES

## TECHNIQUES

- logs, % change, residuals
- point graph [2d histogram], histogram, percentile graph [and with comparisons/reference line], box plot [Tukey]
- dot charts - best way to attach label to quantity, 2-way dot chart {multiway} grouped dot chart
- overlap is dealt with by jitter, distinguishable symbols {sunflower plots}, taking log or other transformation
- box plots for high multiples
- visually distinguish curves and points [this has gotten easy by now]

## CLEVELAND - 3 OR MORE VARIABLES

### THREE OR MORE VARIABLES

- Framed-rectangle graphs
- scatterplot matrices
- interaction/brushing
- 3d wireframe or stereogram (points)

# CLEVELAND - PERCEPTION

## PERCEPTION

- pie v. dot chart
- distance and detection
- length in a stacked bar
- 45 degree banking [Tufte also recommends 1.5:1 horizontal to vertical ratio]
- strive for clarity

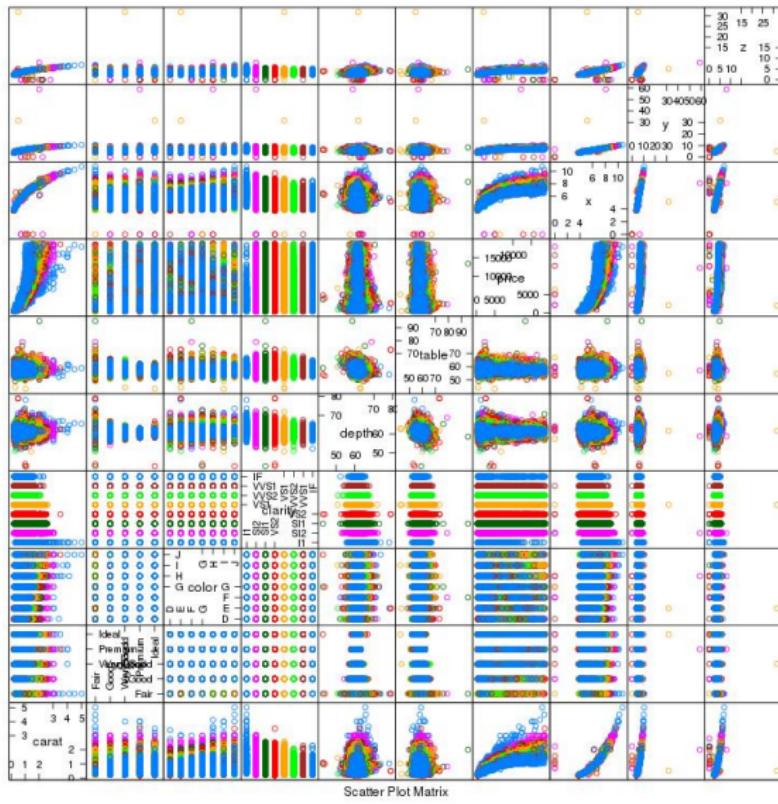
## LATTICE

- The `lattice` package implements Trellis graphics in R
- `lattice` excels at comparative plotting
- uses similar syntax to base graphics, but with greater sorting and manipulative power

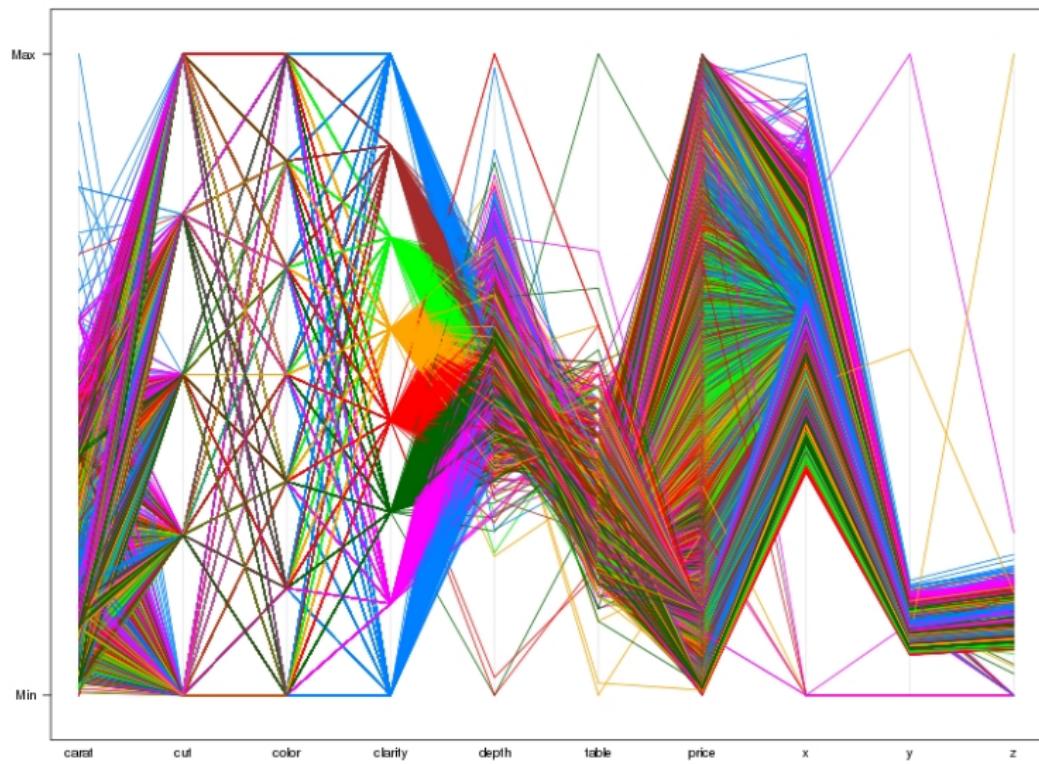
# R BASICS

- `install.packages`, `read.table`, `read.csv`, `library(foreign)`, `library(xlsx)`
- help via `help.start()`, `? package`, `library(help="package")`, `vignette`
- [Cookbook for R](#) and [Quick-R](#) are good jumping off points
- [Guerilla Guide to R](#) or [R crib sheets](#) or [Rutgers R Libguide](#)
- search [Stack Overflow](#) for R graphics to get specific coding tips on graphics
- [R Graph Gallery](#) has many examples of graphs with code, several adapted for this workshop
- [Task Views](#) or [Crantastic](#) can be used to discover packages. In addition to help and online documentation, packages are often described in articles published in places like the *Journal of Statistical Software*.
- This workshop will allow you to use the R language, but glosses over many details of structure and syntax to focus on the graphical elements
- You can always try [R for cats](#) if you get stuck!

# SPLOM



# PARALLEL PLOT



# THE GRAMMAR OF GRAPHICS

*The Grammar of Graphics*, by Leland Wilkinson, was extremely influential in thinking about graphics

- Grammar means "rules for art and science"
- The Grammar of Graphics specifies rules both mathematical and aesthetic
- Earlier graph producers focused on aesthetics of static content
- Dynamic graphics and scientific visualization, by contrast, require sophisticated designs to enable brushing, drill-down, zooming, linking
- The Grammar of Graphics is easily adapted to this approach

# GRAMMAR OF GRAPHICS - PT I

- DATA - weighting, reshaping, counting, bootstrapping
- VARIABLES - transform, sort, log, ranking, residuals, quantiles
- ALGEBRA - nesting or blending data
- SCALES - nominal, ordinal, interval, ratio must be specified
- STATISTICS - static methods available to all graph types e.g, mean, sd, smoothing

# GRAMMAR OF GRAPHICS - PT II

- GEOMETRY - line, area, etc., along with modifiers like jitter and dodge
- COORDINATES - refers to the coordinate system of the graph (cartesian, polar, etc.)
- AESTHETICS - color, texture, size, position, etc. of the data points. Includes using color to classify.
- FACETS - subgroups, multiway tables
- GUIDES - legends, axes, color scales, keys

# GGPLOT2

- `ggplot2` was developed by Hadley Wickham as an implementation of the Grammar of Graphics
- Relatively complete and powerful graphics package
- Can do many things, but not 3D
- See [ggplot2 Help Docs](#) and the `ggplot2` book for complete descriptions
- Other short introductions to `ggplot2` are available, such as these from [Sharp Statistics](#) and [inundata](#)

# A MISCELLANY OF VISUALIZATIONS

- The Cleveland dot chart
  - use to compare labeled quantities, ordered lists
- Kernel Density plot
  - visualize the distribution of data with more precision than a histogram
- Scatterplot Matrix
  - study relationships between all variable combinations

# VISUALIZING DISTRIBUTIONS OF DATA

- Box and Whiskers Plot
  - illustrate quantiles and outliers. There is also a [Tufte version](#).
- Stem and Leaf Plot
  - see precise quantities associated with distribution of data
- Violin plot
  - Blends density information with box and whiskers style (in an artistic manner)
- Dot plot
  - plot distribution point-by-point
- Heat map
  - compare many magnitudes easily

# CATEGORIZING DATA

Many techniques are available to automatically identify related data:

- A *tree* illustrates a categorical classification of the data based on its own characteristics
  - one implementation is the **party** package
- *Self-organizing maps* are a form of neural network that derives characteristics from the data and plots patterns
  - see **kohonen** and **som** packages
- *Clustering* of data can be accomplished by numerous algorithms
- **hclust** and **pvclust** are some methods described at Quick-R's [Cluster Analysis](#) page.
- Other graphs and network analysis tools are available [not explored here]

# VISUALIZING CATEGORICAL DATA

- The *mosaic plot* allows multiple categories to be displayed on the same graph, but can be complicated to interpret.
- The *spineplot* is a variant of the mosaic plot, plotting proportions in 2 dimensions.
- You can also do a [timeline](#).

# MAPS AND GLYPHS

Maps are obviously an important and widespread way of presenting data.

- We examine a few examples of *choropleth* maps, in which shading indicates data levels

*Glyphs* present iconic representations of data elements as plotted points.

- Weather maps often use glyphs.
- A more dynamic example is [here](#).
- As an R example, consider Chernoff faces and the `aplypack` package. Also, [Smiley faces](#) [and many more graph variants in this chapter].

# 3-D

- 3-D scatterplots
  - `cloud` (`lattice`)
- contour plots
  - to plot standardized levels of data
- wireframe plots
  - to present a 3-D surface representation of data
- `rgl` (a separate package containing several 3d plotting functions and animation)
- `mosaic3d` extends the mosaic paradigm to three dimensions

# ANIMATION

- Animation is an easy way to step through data over time
- or to provide comparisons of different views of data
- R makes animation easy with the `animation` package
- Just enclose a sequence of graphics in the `animation` command to generate interactive HTML (or GIF, SWF, L<sup>A</sup>T<sub>E</sub>X, Video).

# INTERACTIVE DATAVIZ - PRINCIPLES

- Why aren't all of our graphs interactive?
- *Brushing* is used to select data points and track them through various analyses.
- Drilling down, zooming, and subsetting are also interactive techniques.
- Data displays can be linked so that a selection in one panel modifies the output displayed in another panel.
- Interactivity is especially useful for data exploration, studying multidimensional relationships.

# LINKED DATA PANELS AND VIS PACKAGES

In many contexts, visualizing the relationships between data elements is made easier by viewing related data panels simultaneously.

- One example of this occurs in time series data with decomposition into trend, seasonal, and random components
- The tableplot (`tabplot` package) implements another linked data view across all variables.
- googleVis and other “Vis” packages, e.g. `bdvis` for biodiversity or `rainfreq`.

# INTERACTIVE DATA IN PRACTICE

There are many R packages that allow for interactive data work in a graphical user interface, including:

- **playwith** - versatile package that works with any graphics function. Graphics can be explored, edited, and exported.
  - requires separate installation of GTK+ on your computer [method varies by OS]
- **rggobi** - powerful 3-D tool for brushing, identifying and manipulating data with a book and online [companion site](#).
  - requires separate installation of GGobi on your computer [method varies by OS]
  - growing long in the tooth [no update since 2012, problematic install on Mac OS X]
- **rattle** - package designed for data mining, includes graphics options alongside other statistical functions
- **latticeist** - allows complex linking of plots

# INTERACTIVE DATA ON THE WEB - RCHARTS

- **Rcharts** is a package that uses javascript to create interactive visualizations.
- Lattice-style commands are used.
- The package can output javascript for use in an HTML page.
- Some commands depend on supplemental javascript libraries that must be installed, such as **NVD3**
- Can embed in documents too, with **slidify**

# INTERACTIVE DATA ON THE WEB - SHINY

- The [shiny](#) package is developed by the Rstudio folks
- You can learn shiny in half a day via the online tutorial
- More custom control of the design is possible with shiny, in comparison to other do-it-all packages
- Graphics use familiar R syntax, with wrappers to implement web functionality
- Every shiny app has the same structure: two R scripts saved together in a directory [ui and server files]
- You must install the shiny server to deliver pages via the web
- For now, we'll just demo a few examples. You can see more in the [shiny gallery](#)
- **Rcharts** works with **shiny** too.

# INTERACTIVE DATA ON THE WEB - GGVIS

- The `ggvis` package is ALSO developed by the Rstudio folks
- Think `ggplot` meets `shiny`
- Similar syntax to `ggplot`
- Some ability to add interactive controls

# BIG DATA

- Big data presents special issues for data visualization
- While many techniques and graphics are the same, exploration and plotting must be optimized for the size of the data set
- Representation of the complexity of the data may require special techniques
- `hexbin`
- `bigvis`

# AIRLINE DATA

For this exercise, we will use [Airline on-time performance data](#).

- This data contains information on every flight in the United States from 1987 to 2008, including arrival and departure times, delays, and other attributes.
- The link above allows access to the full dataset. These are large. For example, the full 2008 data contains 7,009,728 records and is 689 MB.
- We will use an extract of selected variables from January 2008 only. This subset contains 605,765 records and is 37.5 MB.

# BINNING

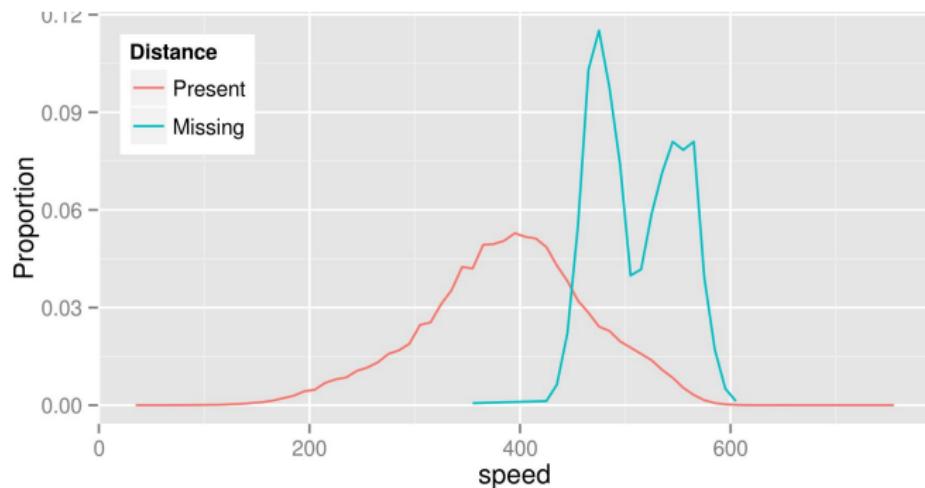
- `hexbin` resolves overplotting issues in large data sets by showing density
- There are many other binning methods

**BIGVIS** is a relatively new package by Hadley Wickham to deal with the issues of Big Data

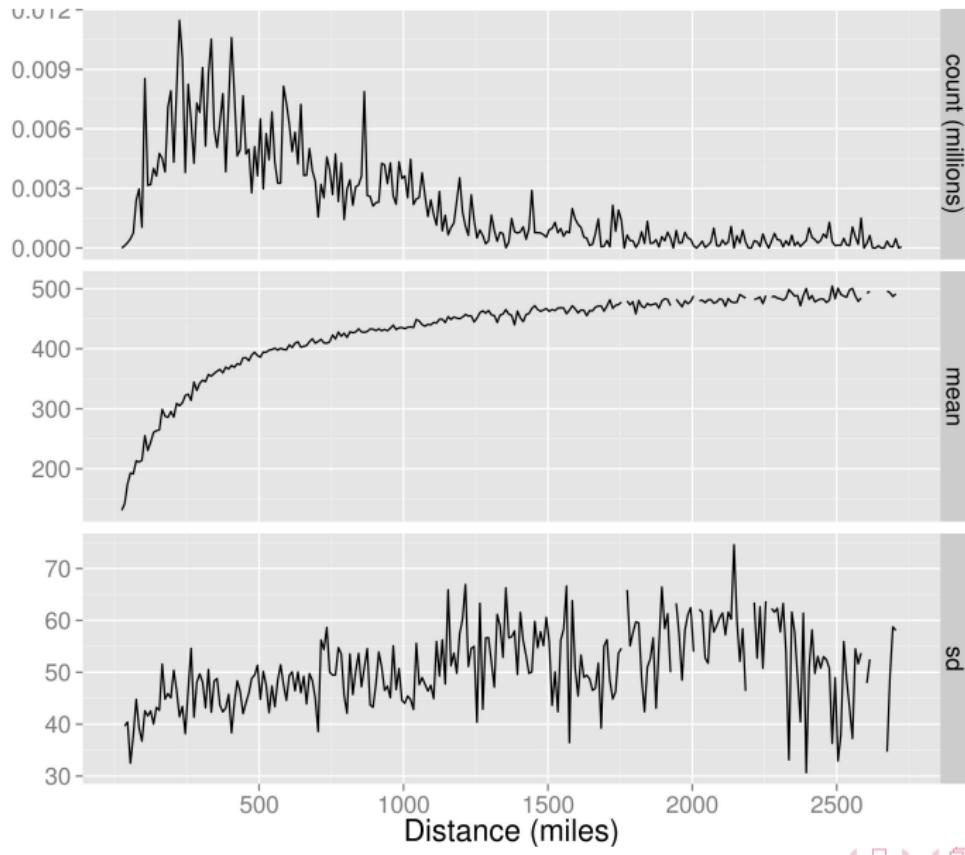
- There is a [Preprint](#) and [R Meetup](#) presentation by Hadley Wickham
- Complete code, including the extracts adapted for this workshop, is available at <https://github.com/hadley/bigvis-infovis>
- Target: process 100 million observations in under 5 seconds.
- Fundamental principle: No need for more data points than there are pixels on the screen.

## BIGVIS STEPS

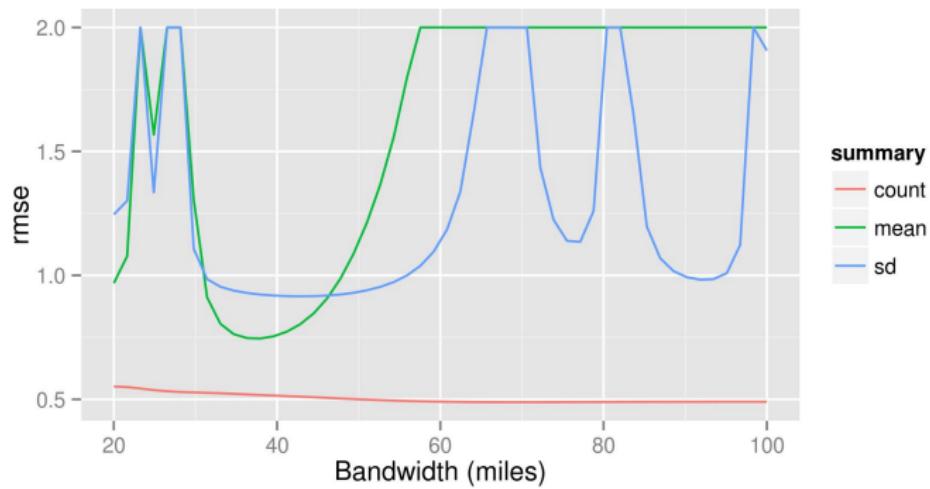
- Condense (`bin, condense`)
- Smooth (`smooth, best_h, peel`)
- Visualize (`autoplot` plus standard methods)
- The following slides are summary output from Hadley's bigvis example.



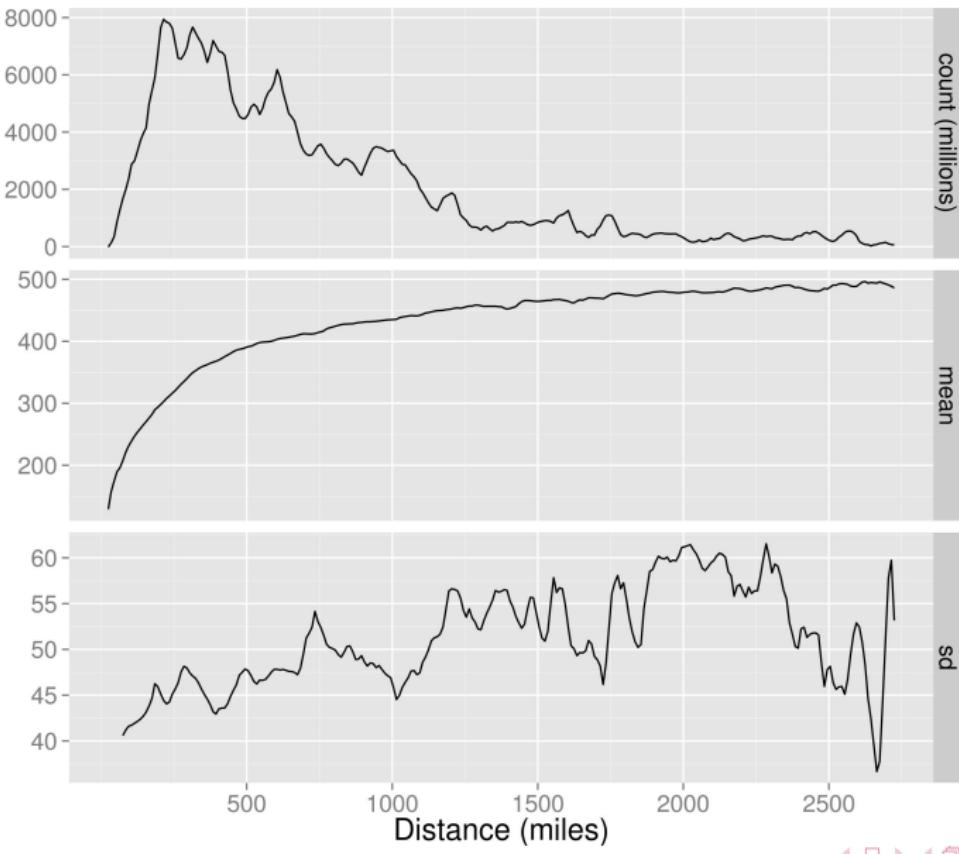
# BIGVIS

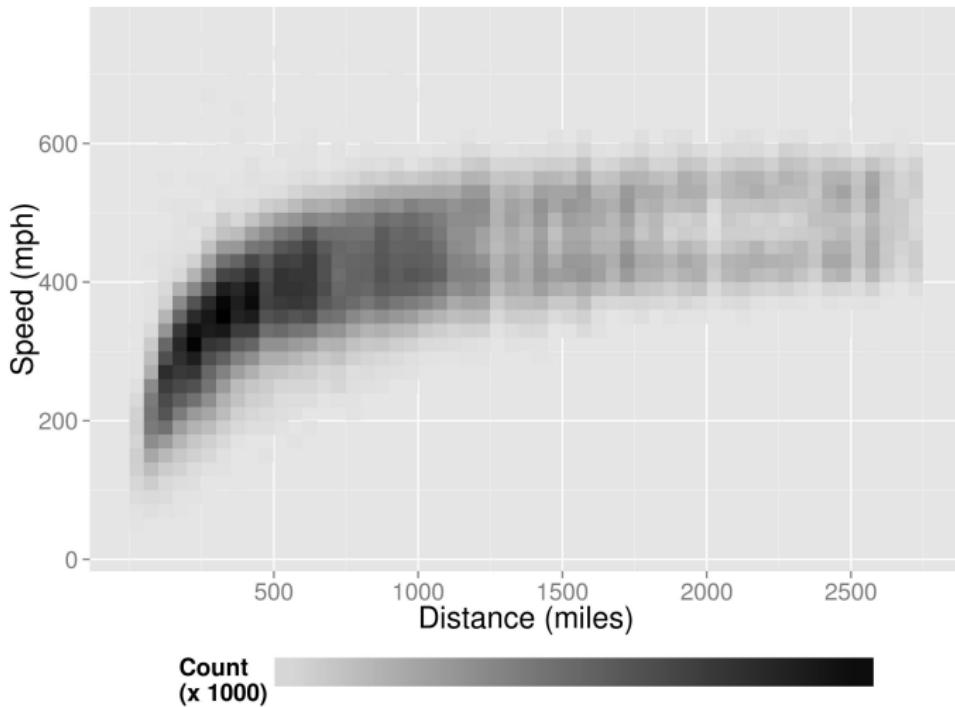


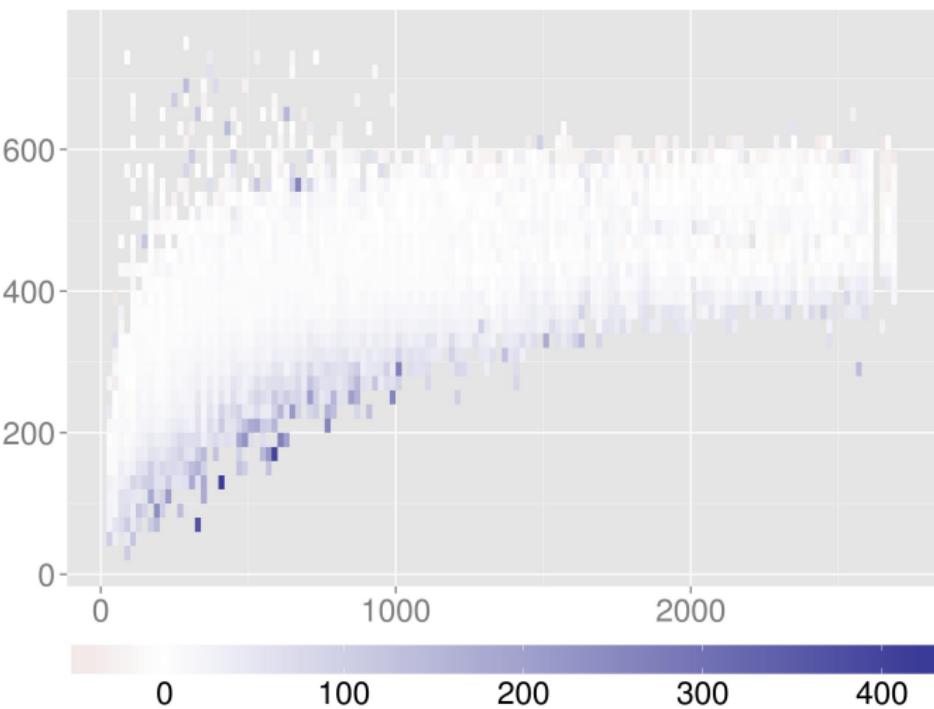
# BIGVIS

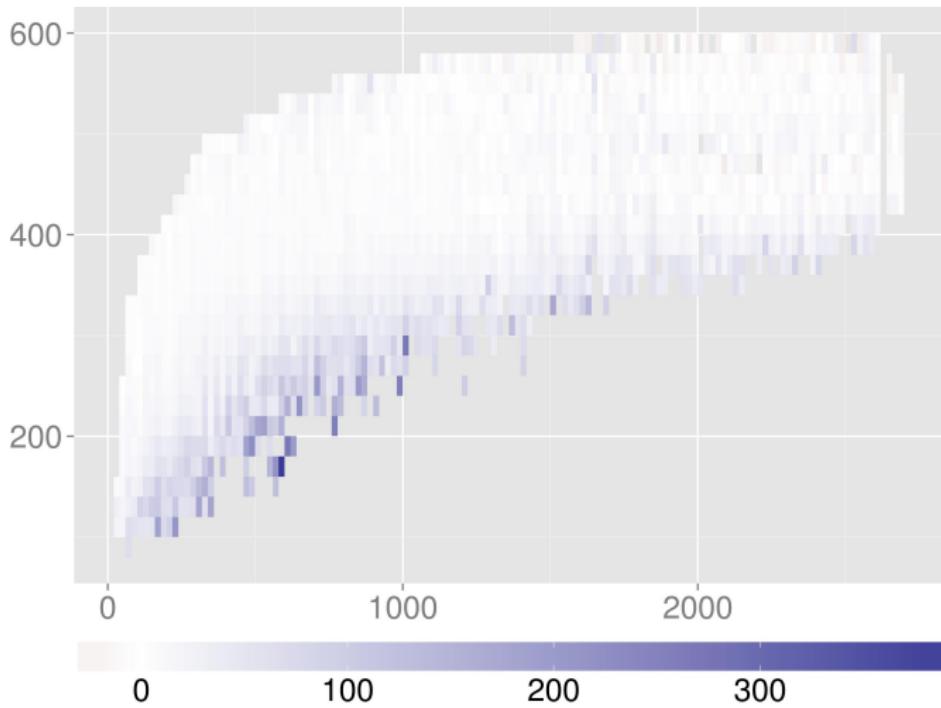


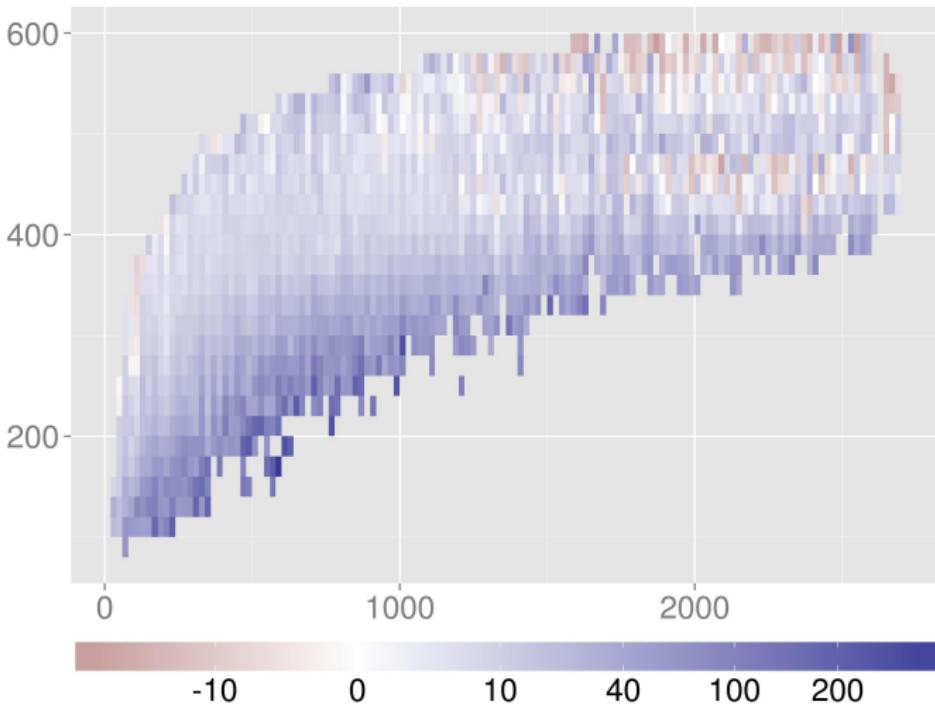
# BIGVIS











# INFOGRAPHICS LINKS

Although not covered here, the following links are a sampling of infographics sites for your later enjoyment:

- Data Storytelling in Video
- Art of Data Visualization - in spite of its title, more on the infographics side
- Parisian Subway Traffic and New York Subway Inequality
- Tulp Interactive
- Information Aesthetics
- Mapping London and London Riots + Twitter
- YouTube Trends Map
- Global Burden of Disease Visualizations

# OTHER CONSIDERATIONS

These are not illustrated in the code, but represent future topics for exploration.

- confidence intervals
- missing data [discuss]
- color can be used to indicate certainty
- “scagnostics”, borderlining scatterplots (**scagnostics** package)
- edge blur, crisp vs. fuzzy edges

# KEEP EXPLORING

Data Visualization represents a nearly infinite world of possibility for exploration:

- plunge into programming
- deep dives into data
- indulge in interactivity
- ...have fun and keep learning! [e.g., [R-bloggers.com](#)]

# REFERENCES

There is also an online bibliography of references to accompany this presentation on [my home page](#).