

Data Science in the Wild

Lecture 13: Information Visualization

Eran Toch



**CORNELL
TECH**

Agenda

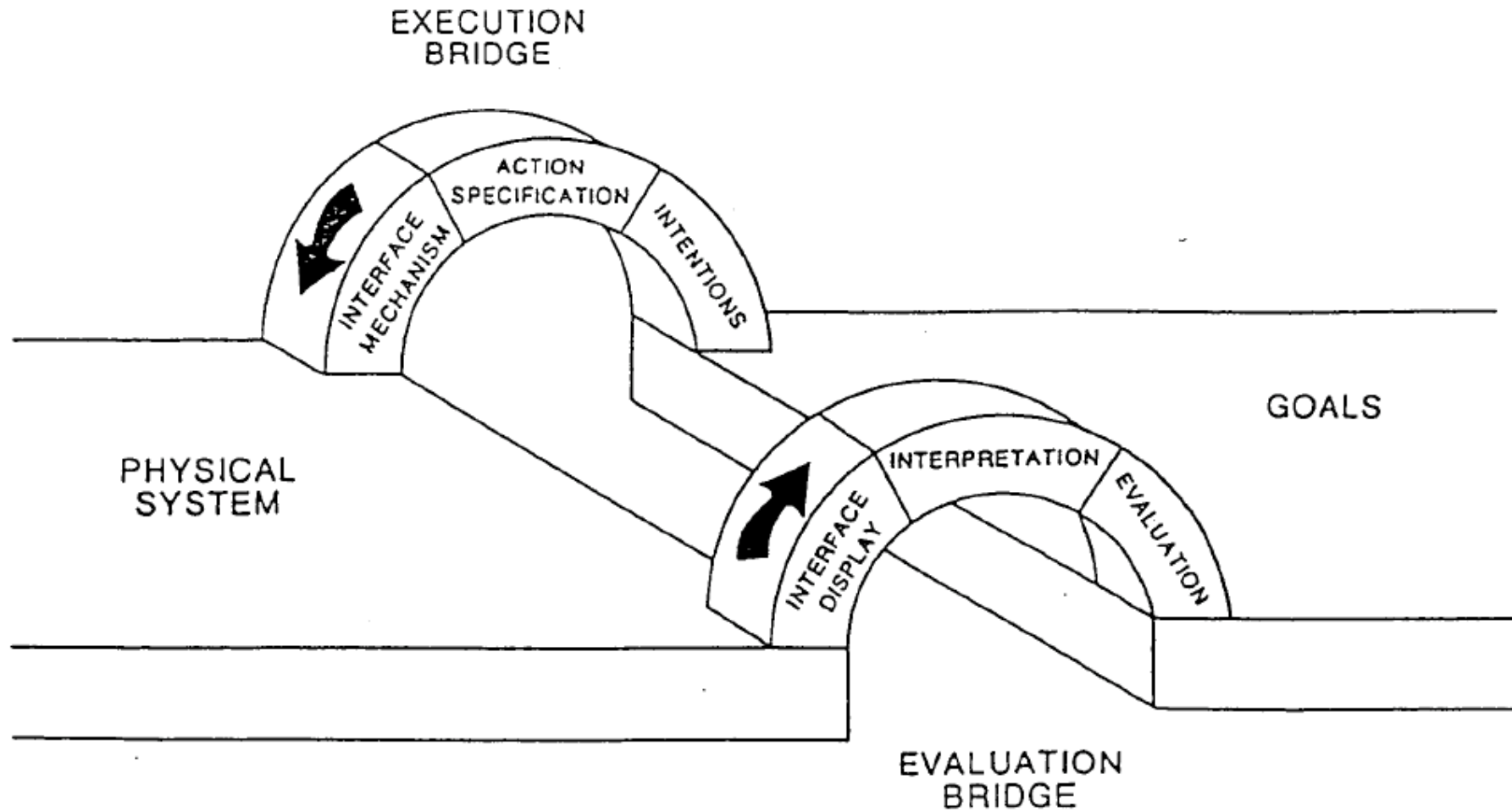
1. What is visualization?
2. Types of visualization
3. Good visualizations

A definition

- “Transformation of the symbolic into the geometric” (McCormick et al., 1987)
- The depiction of information using spatial or graphical representations, to facilitate comparison, pattern recognition, change detection, and other cognitive skills by making use of the visual system

By Marti Hearst

Gulf of Execution and Evaluation

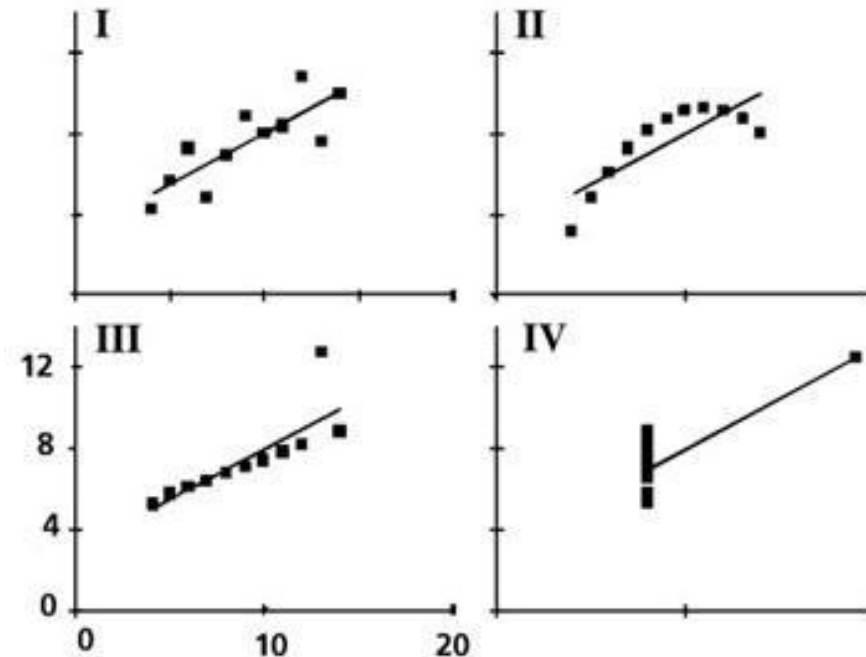


Norman 1986

The limitations of symbolic data

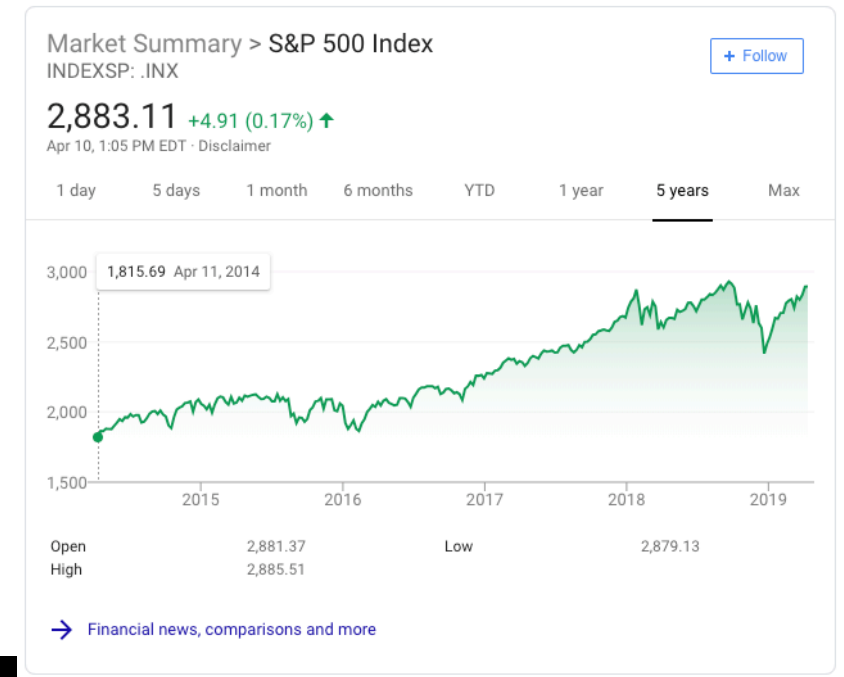
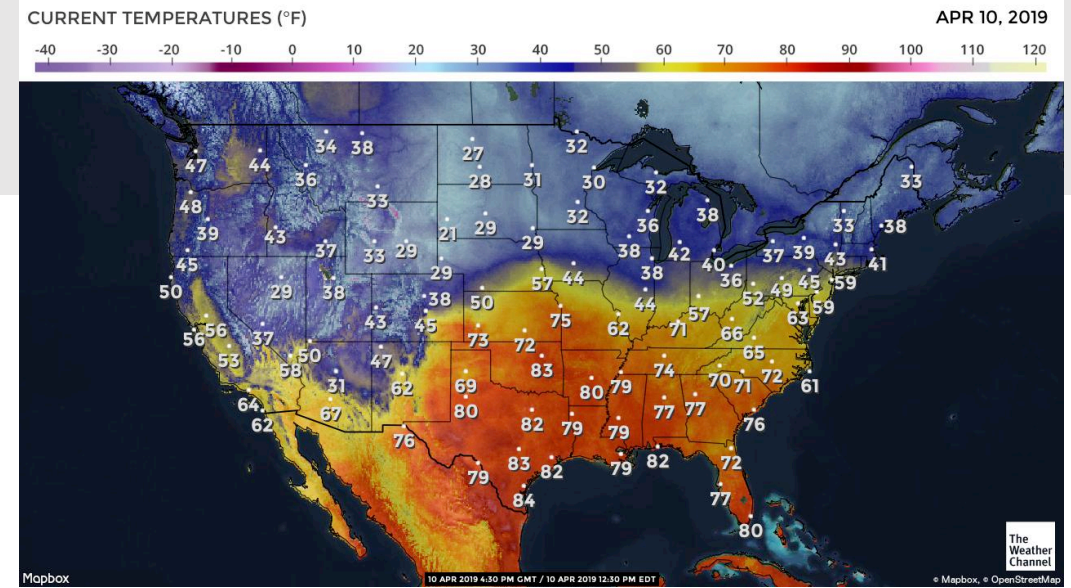
The regression lines for all these data points are exactly the same

| I | | II | | III | | IV | |
|------|-------|------|------|------|-------|------|-------|
| x | y | x | y | x | y | x | y |
| 10.0 | 8.04 | 10.0 | 9.14 | 10.0 | 7.46 | 8.0 | 6.58 |
| 8.0 | 6.95 | 8.0 | 8.14 | 8.0 | 6.77 | 8.0 | 5.76 |
| 13.0 | 7.58 | 13.0 | 8.74 | 13.0 | 12.74 | 8.0 | 7.71 |
| 9.0 | 8.81 | 9.0 | 8.77 | 9.0 | 7.11 | 8.0 | 8.84 |
| 11.0 | 8.33 | 11.0 | 9.26 | 11.0 | 7.81 | 8.0 | 8.47 |
| 14.0 | 9.96 | 14.0 | 8.10 | 14.0 | 8.84 | 8.0 | 7.04 |
| 6.0 | 7.24 | 6.0 | 6.13 | 6.0 | 6.08 | 8.0 | 5.25 |
| 4.0 | 4.26 | 4.0 | 3.10 | 4.0 | 5.39 | 19.0 | 12.50 |
| 12.0 | 10.84 | 12.0 | 9.13 | 12.0 | 8.15 | 8.0 | 5.56 |
| 7.0 | 4.82 | 7.0 | 7.26 | 7.0 | 6.42 | 8.0 | 7.91 |
| 5.0 | 5.68 | 5.0 | 4.74 | 5.0 | 5.73 | 8.0 | 6.89 |

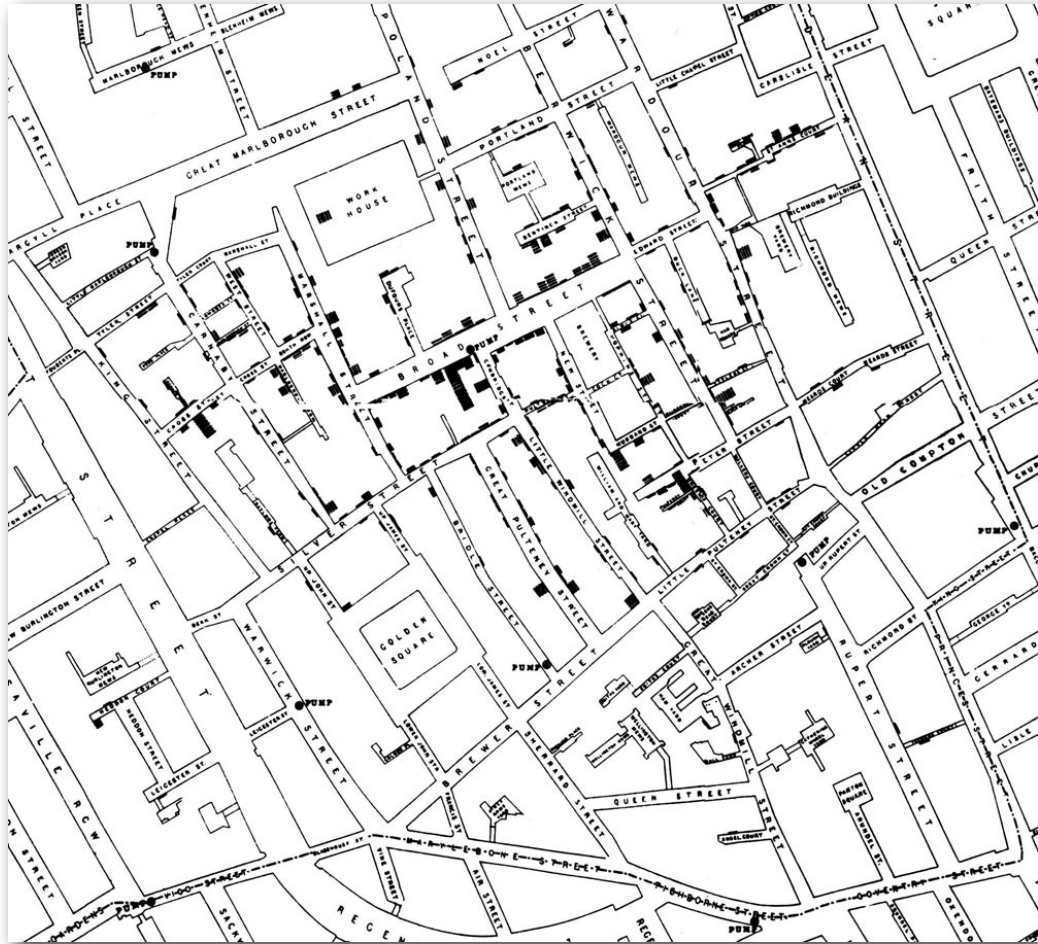


When do we use visualization

- **Explore the data**
 - Make large datasets accessible
 - Support scanning, outlier detection, recognizing and exploring
- **Support analysis**
 - Pattern recognition
 - Comparisons
- **Tell stories about the data**
 - Explain
 - Reason

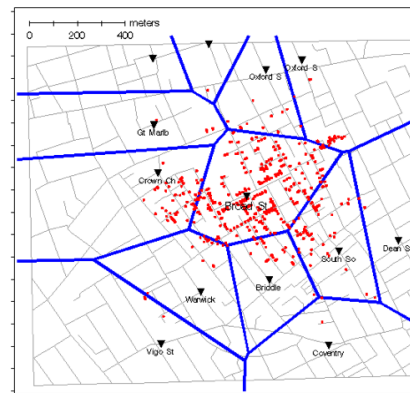


Good Visualizations



Dr. John Snow's map of deaths from a cholera outbreak in London, 1854, in relation to the locations of public water pumps.

London cholera deaths, 1854: polygons

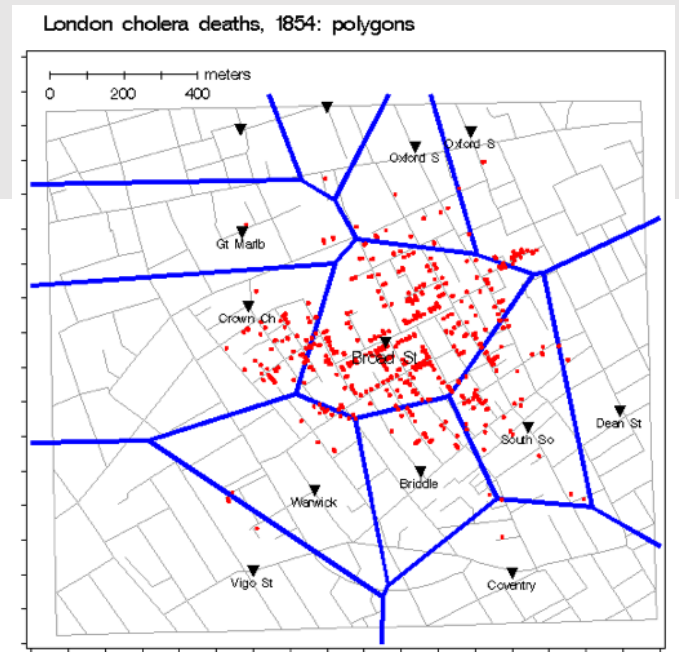


Broadwick Street showing the John Snow memorial and public house.

Theory of Visualization

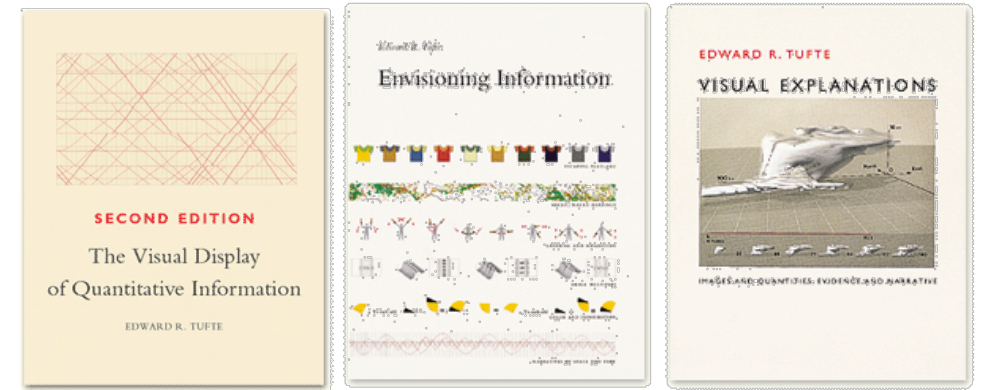
Mapping data questions to proper visualization.

- How to display dimensions of different types
- How to use color? Shapes?
- How to present multiple dimensions?
- What are the cognitive properties of visualization?
- How to differentiate between good and bad visualization?



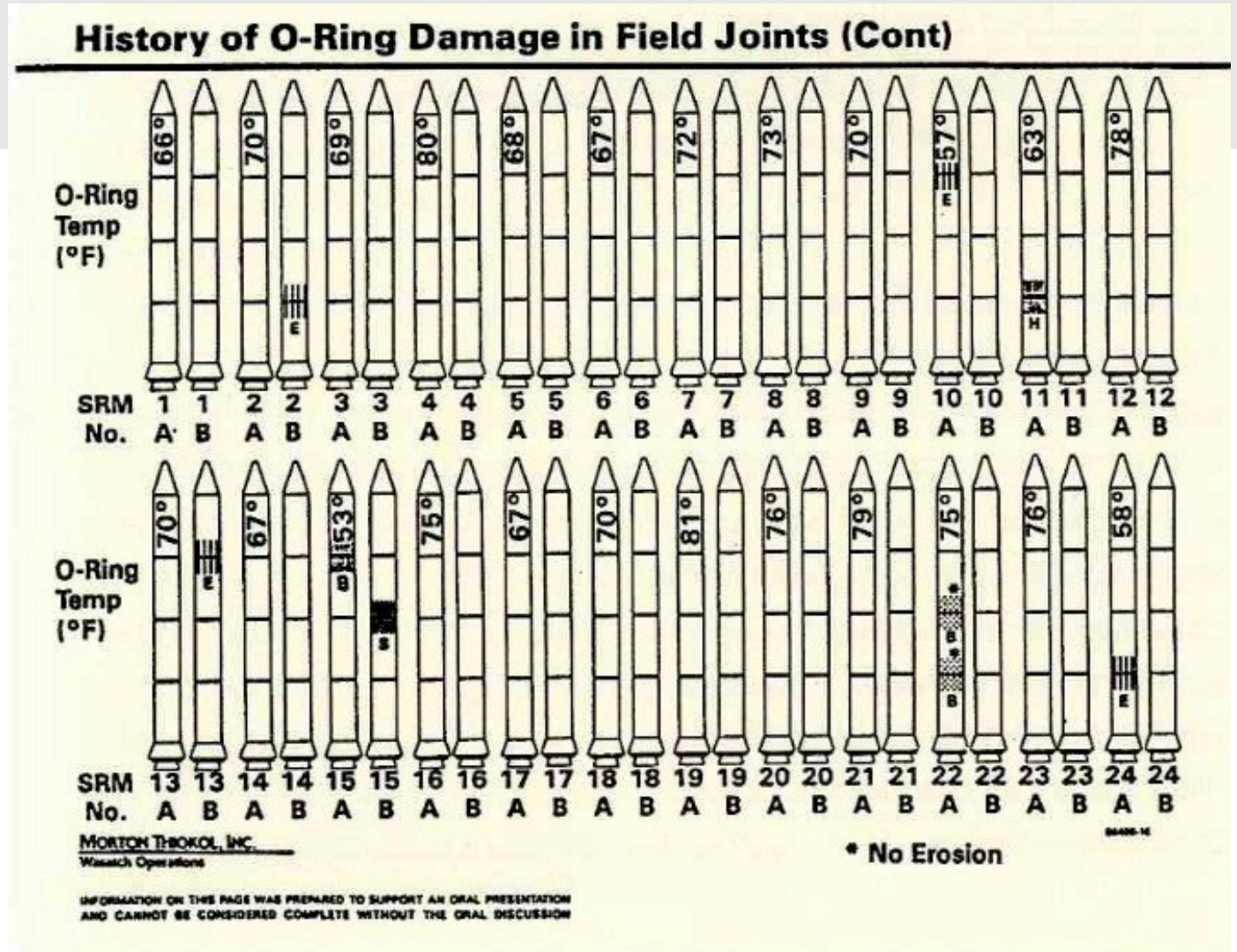
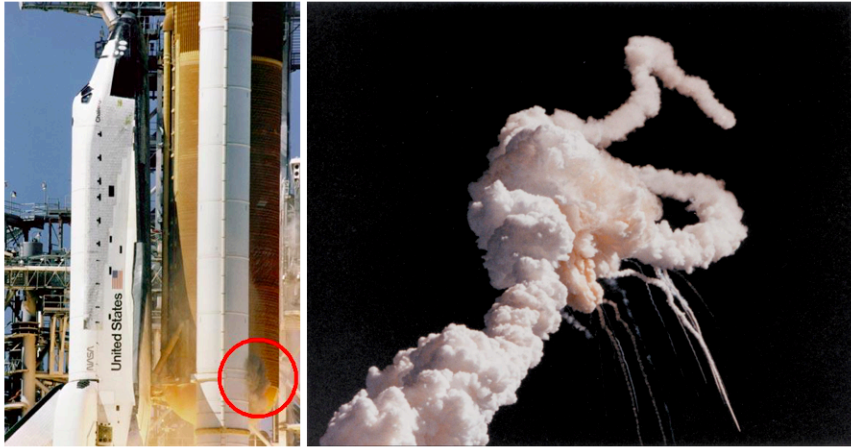
Theory

- Kosslyn: Types of Visual Representations
- Lohse et al: How do people perceive common graphic displays
- Bertin, MacKinlay: Perceptual properties and visual features
- Tufte/Wainer: Representing honesty and context



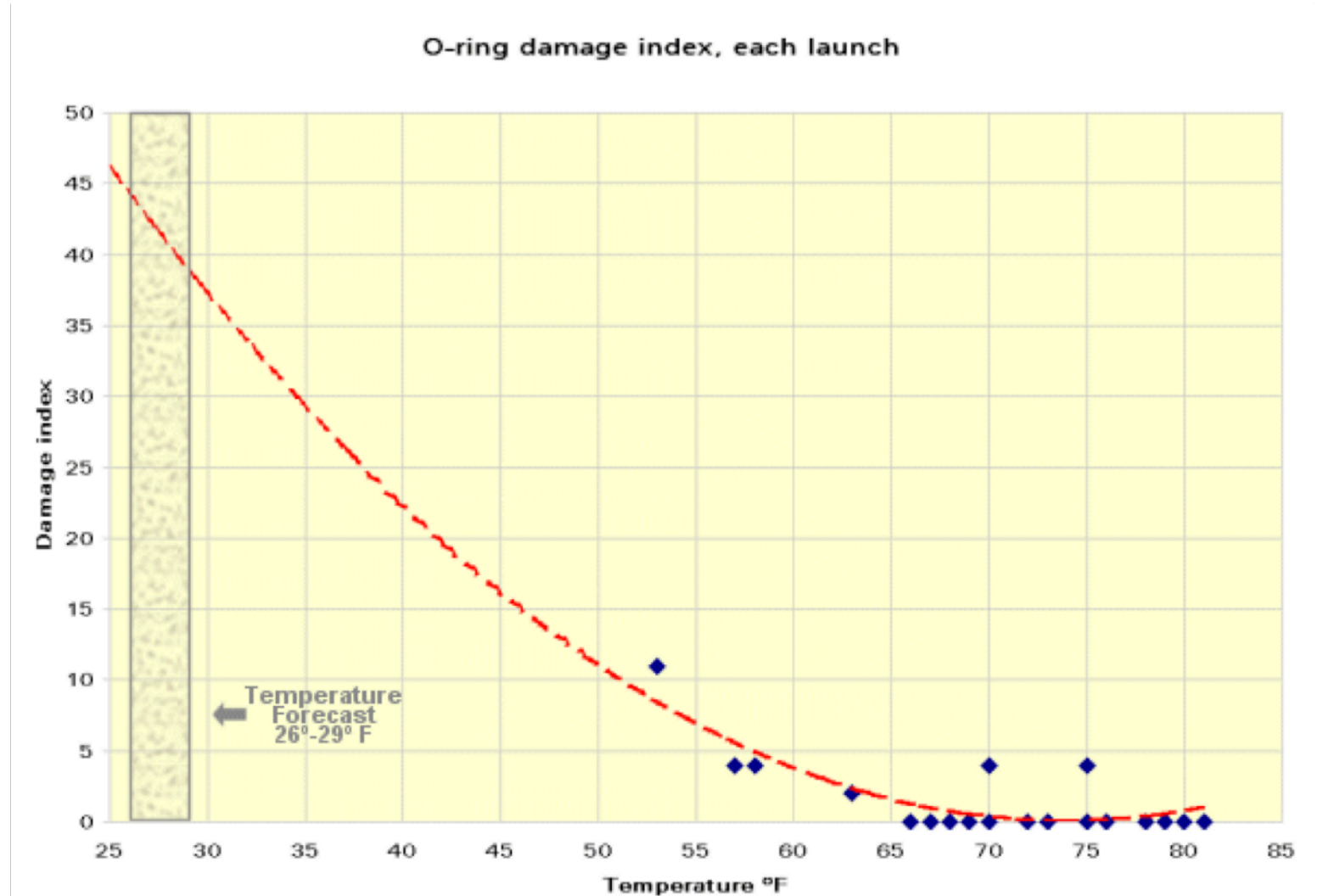
Tufte

Bad Visualization



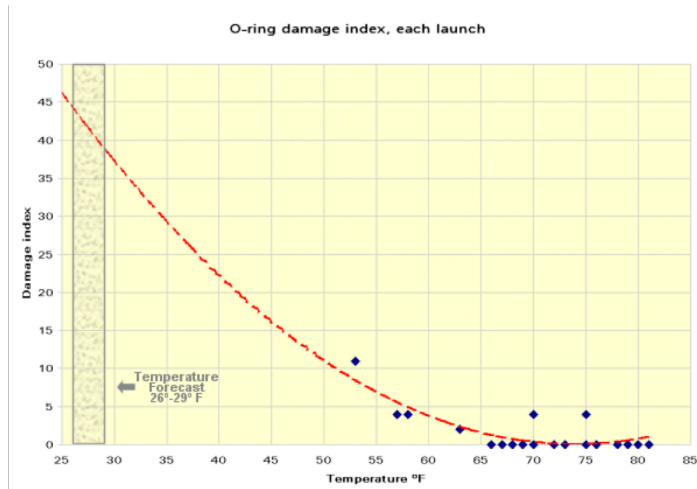
The damage to the O-rings and the temperature is associated with each launch

A (possibly) better visualization



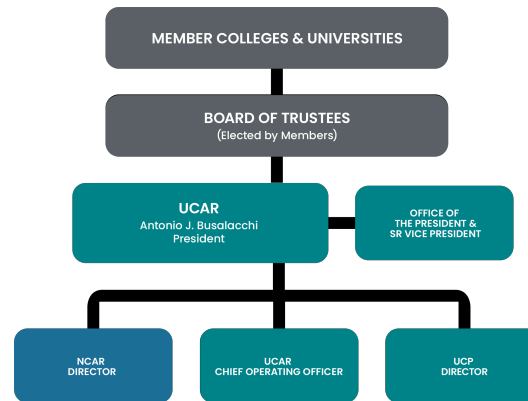
<https://medium.com/the-data-experience/four-questions-you-should-ask-before-visualizing-your-data-cd20a302eb65>

Visualization Types



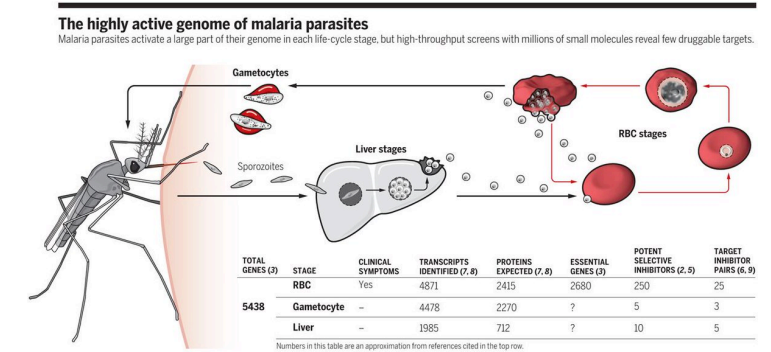
Graph

At least two scales, one of them quantitative



Chart

Discrete relations and structures between entities

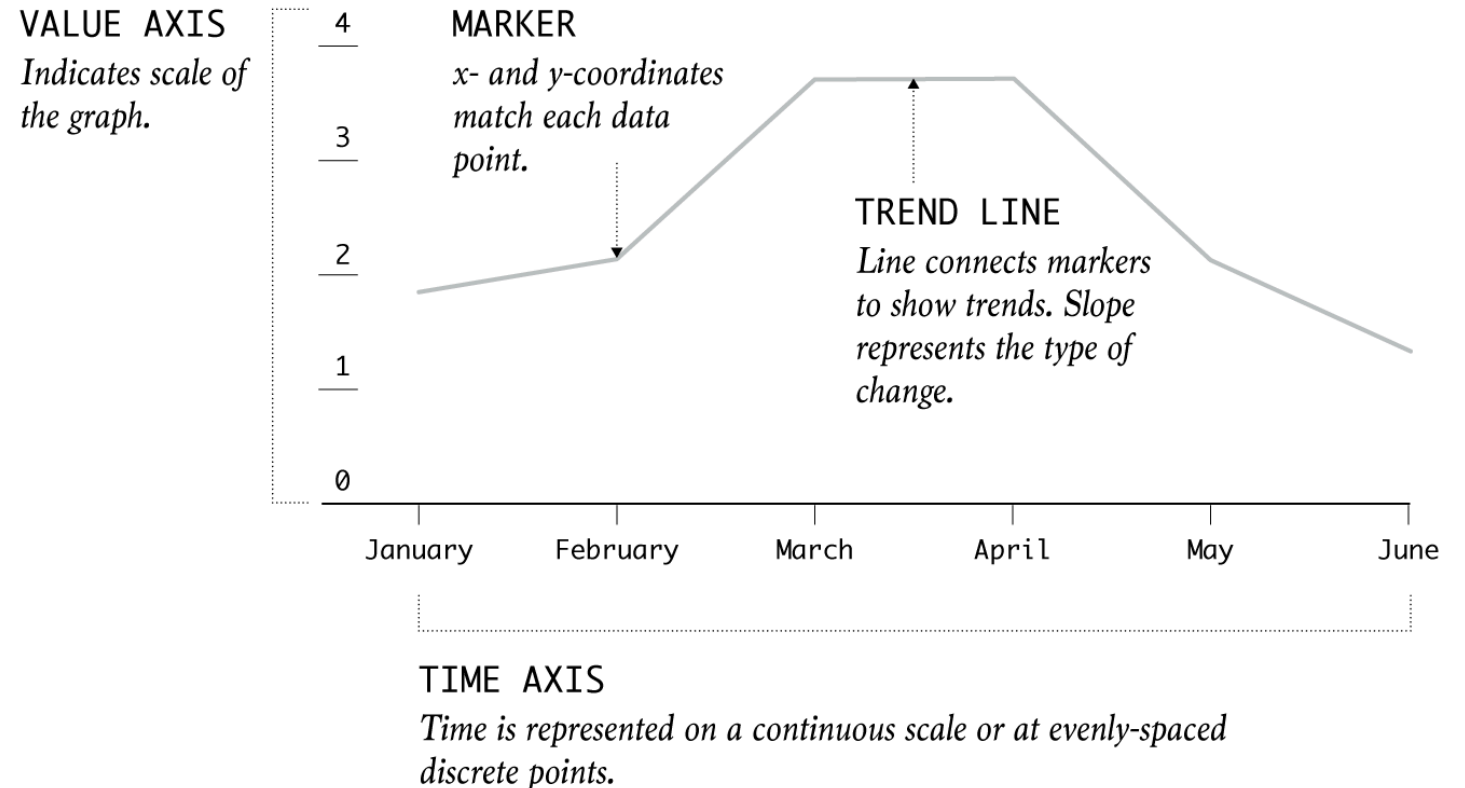


Diagram

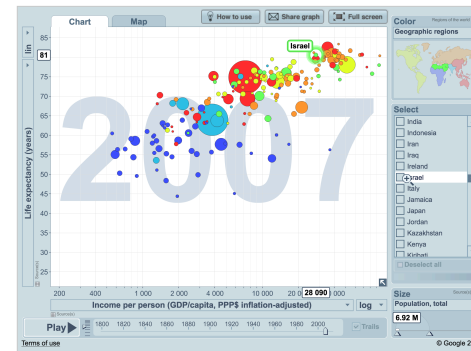
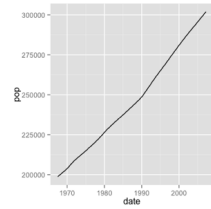
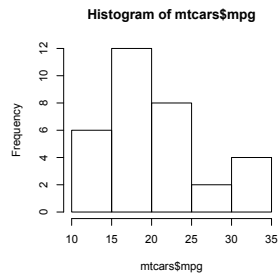
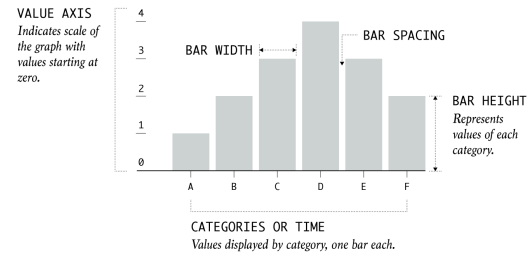
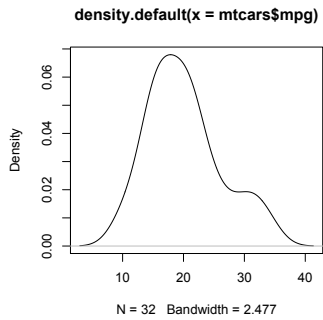
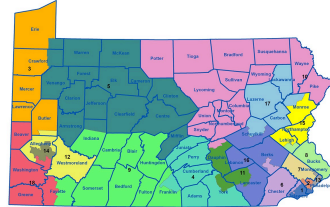
Mixing graphic schematics and symbolic information

What makes a graph?

- Framework
 - Scale
 - Layers
- Content
 - Graphics: points, lines, areas, bars
 - Marks
- Labels
 - Titles
 - Axes
 - Tic marks

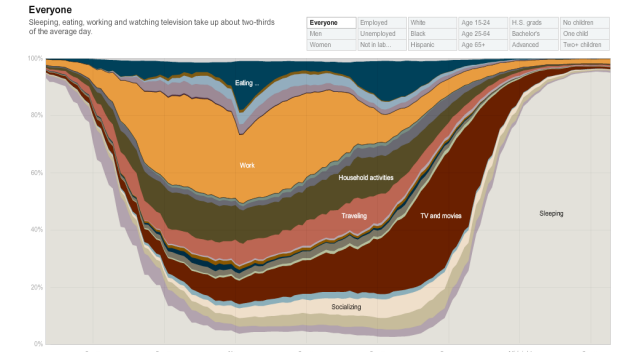


Graph types

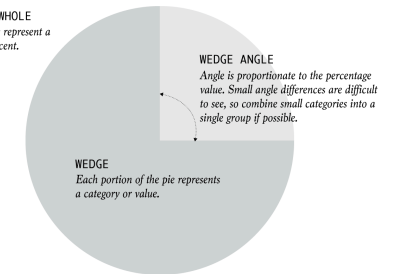


How Different Groups Spend Their Day

The American Time Use Survey asks thousands of American residents to recall every minute of a day. Here is how people over age 15 spent their time in 2008. Related article



PARTS OF A WHOLE
Sum of all wedges represent a whole, or 100 percent.



Distribution

Comparison

Relationship

Composition

Accuracy and Visualization

Less Accurate

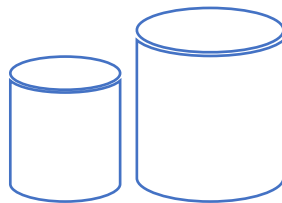
More Accurate



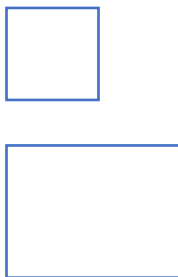
Color



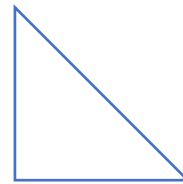
Volume



Area



Slope



Length



Position



Density

Angle

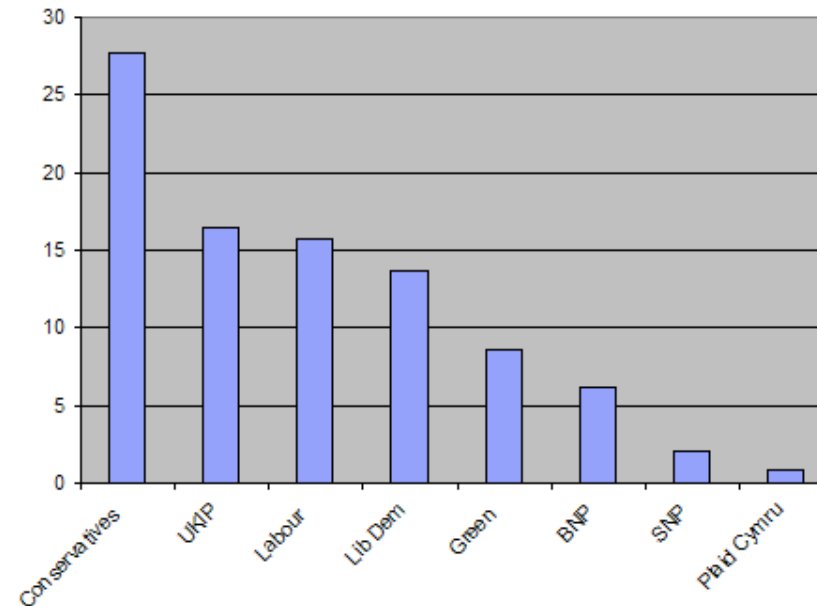
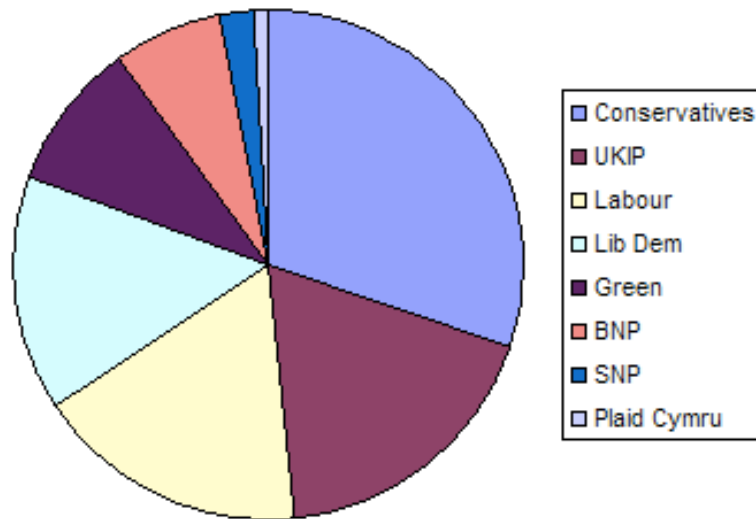
[Mackinlay 88 from Cleveland & McGill

World's Most Accurate Pie Chart



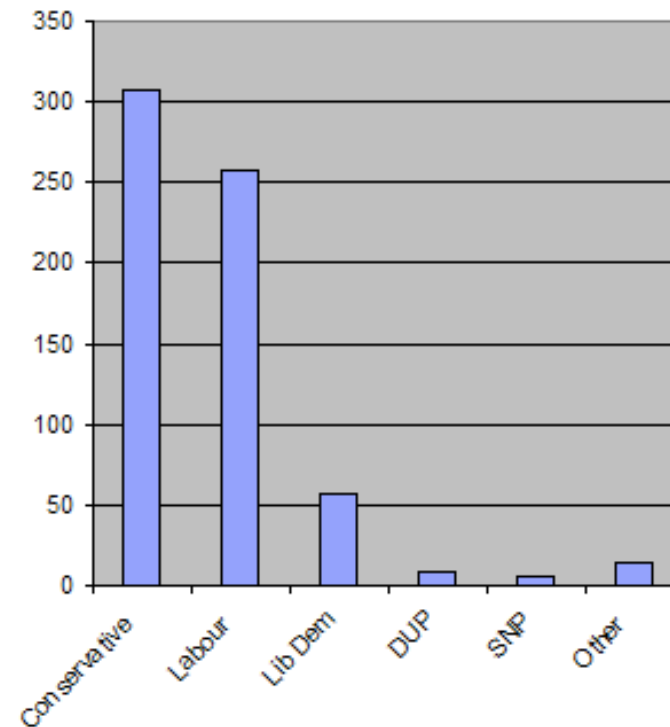
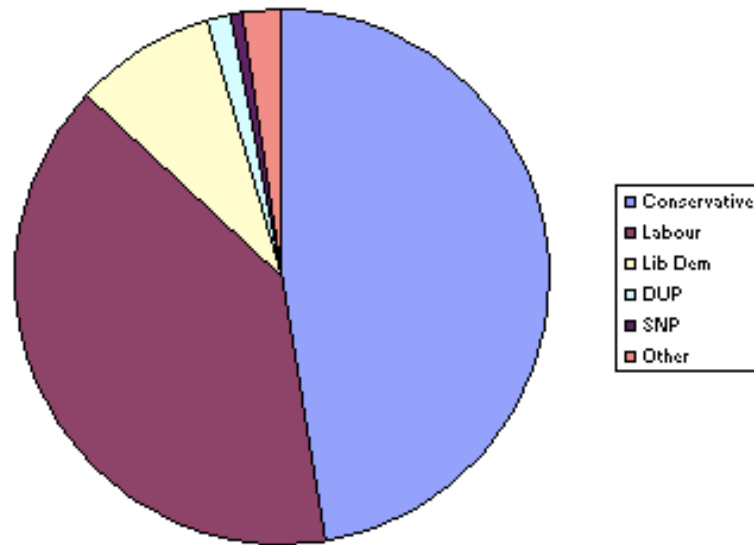
Hey, what do you have against Pie Charts?

- Hard to compare similar values
 - With more than two values
- What is the difference between the second, third and fourth parties?



Pie-Charts are Evil

- Difficult to evaluate combined values.
- Can the Labour and Lib-Dem parties form a coalition?



Good Visualization Theory

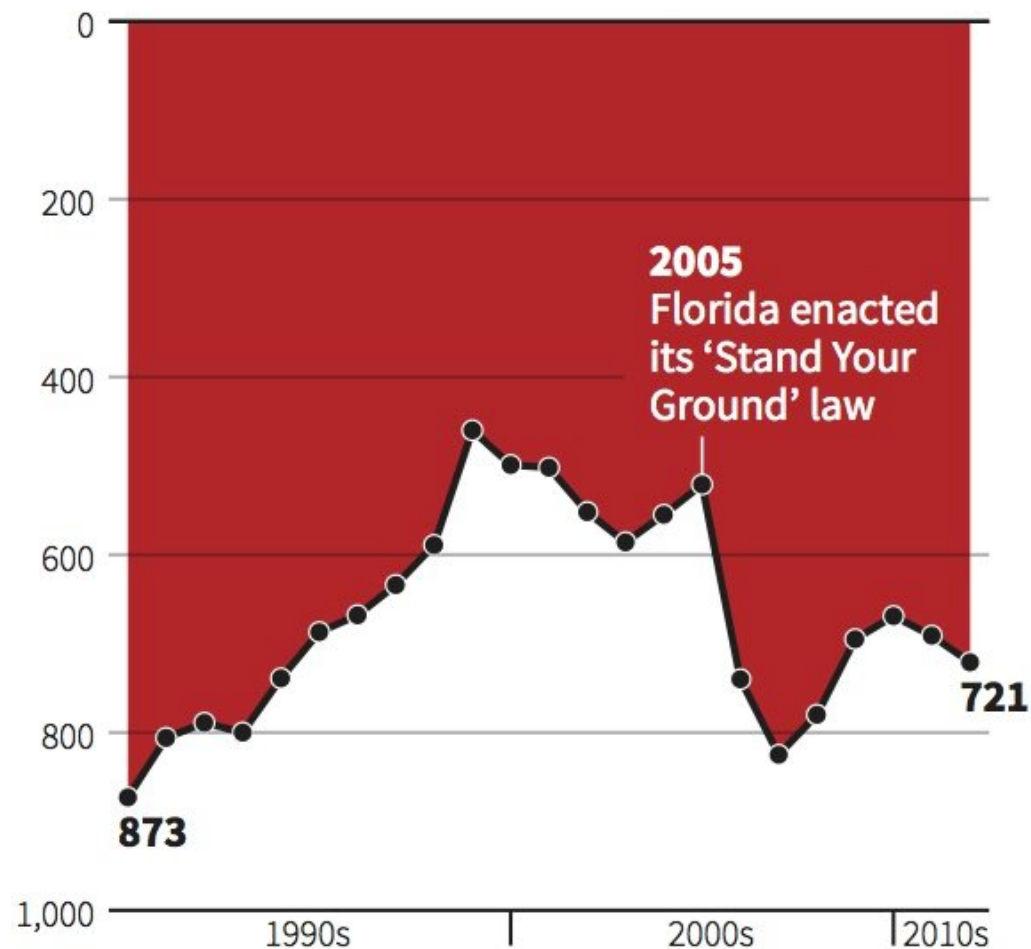
- Tufte defines several rules for creating good visualizations:
 - Avoid distorting the data
 - Maximize the data-ink ratio
 - Every pixel requires a reason
 - Encourage Eye to Compare Different Pieces of Data
 - Putting data in context: linear averages, examples
 - Closely integrate Statistical and Verbal Descriptions

Data Distortion

- Graphical integrity:
 - Misleading uses of area
 - Misleading uses of perspective
 - Leaving out important context

Gun deaths in Florida

Number of murders committed using firearms



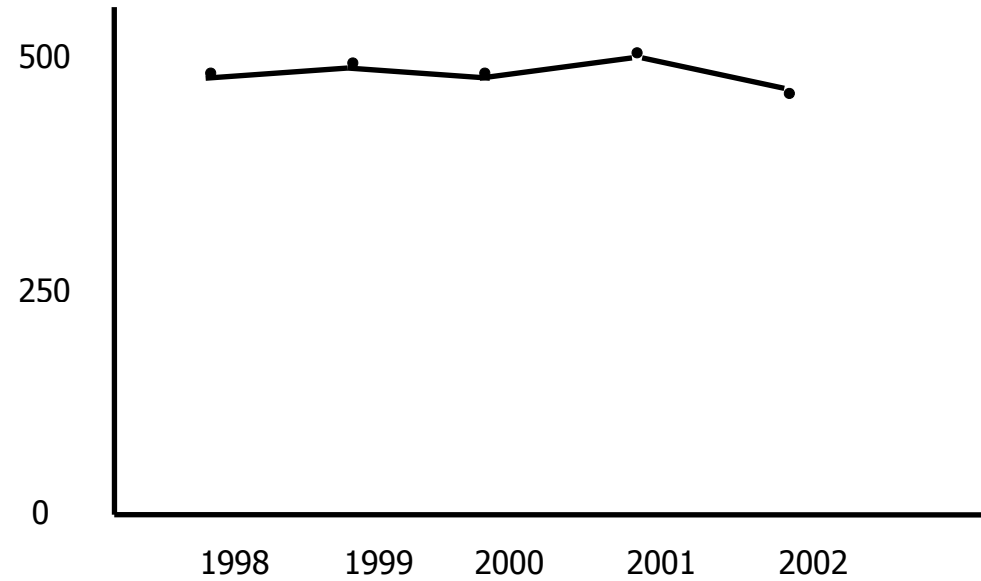
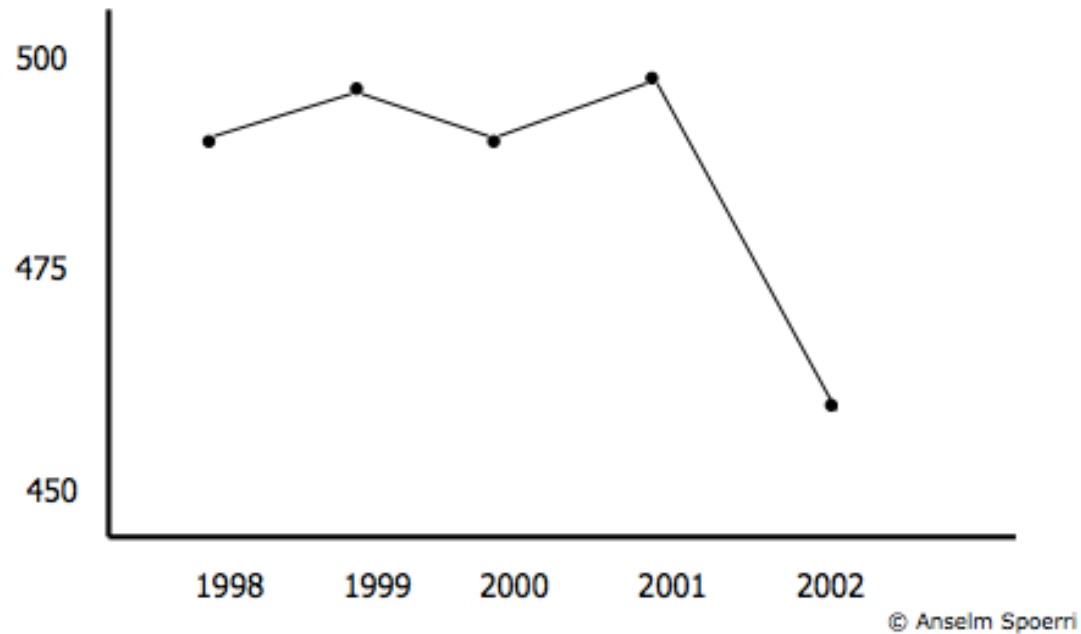
Source: Florida Department of Law Enforcement

C. Chan 16/02/2014

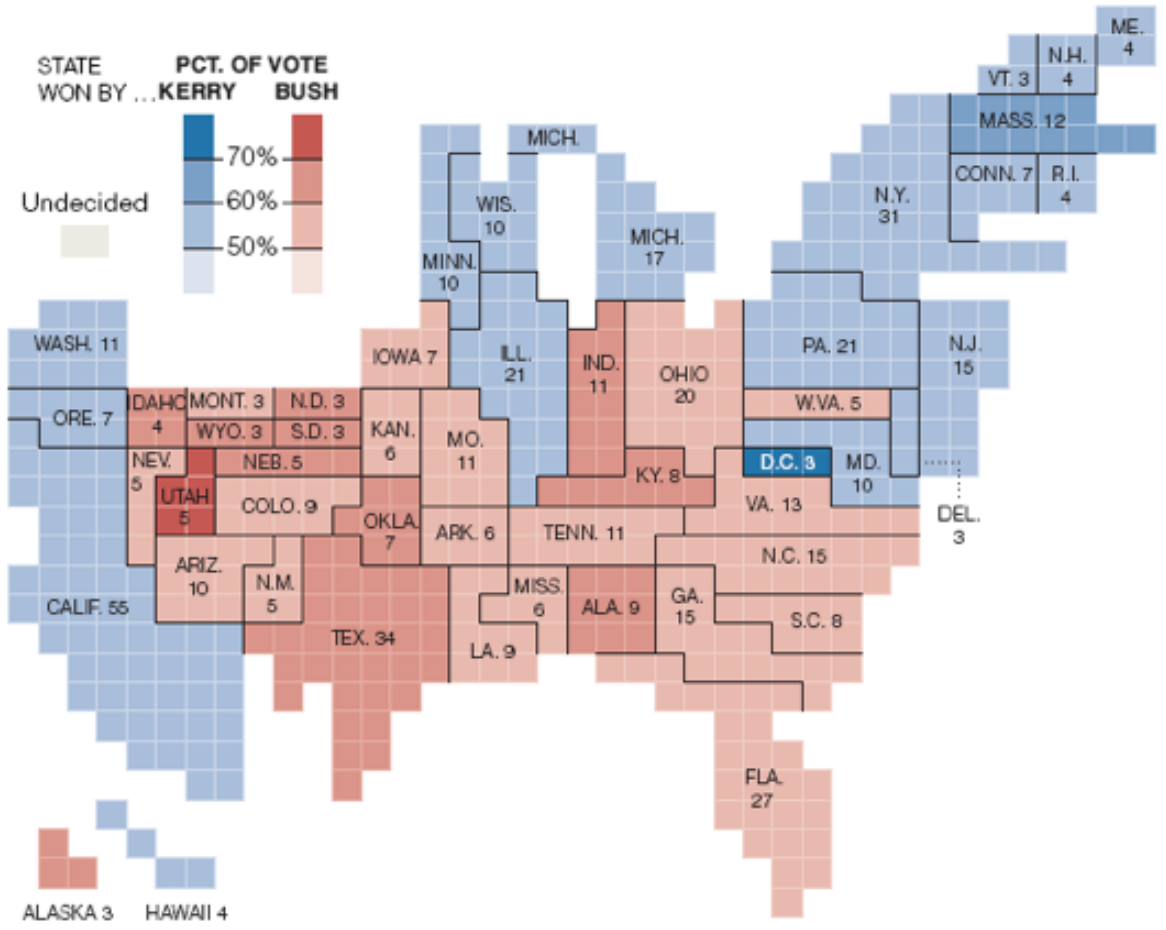
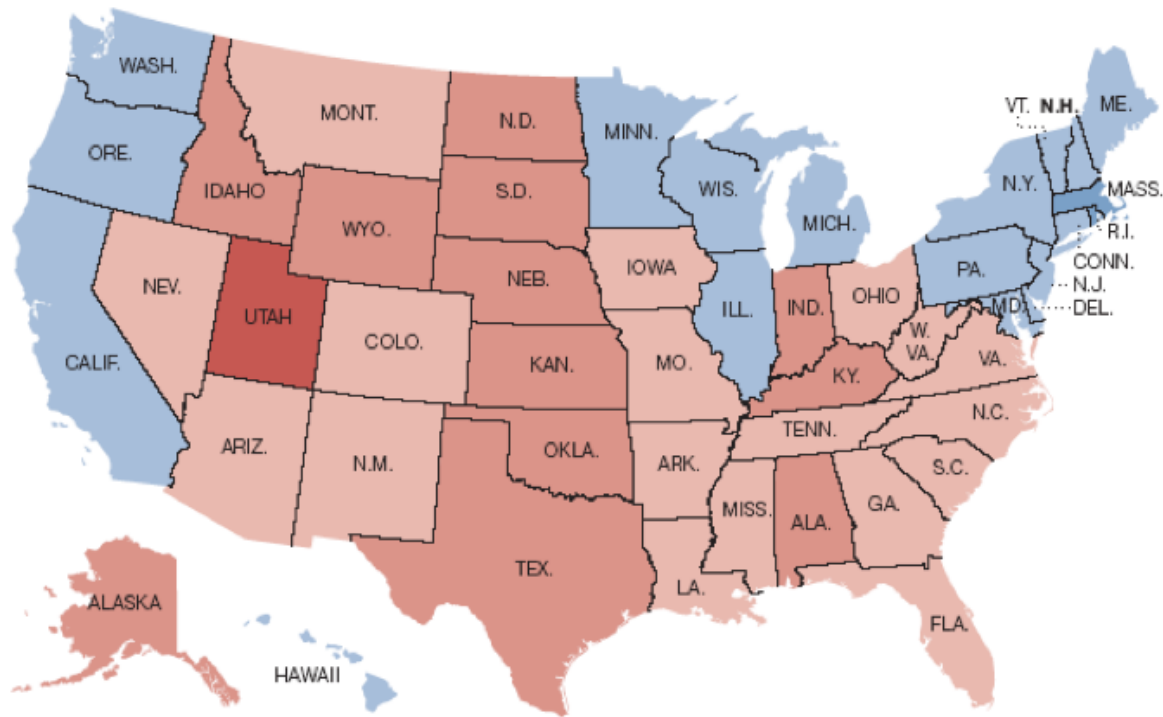
REUTERS

Scale

Lie factor = $\frac{\text{size of effect in graph}}{\text{Size of effect in data}}$

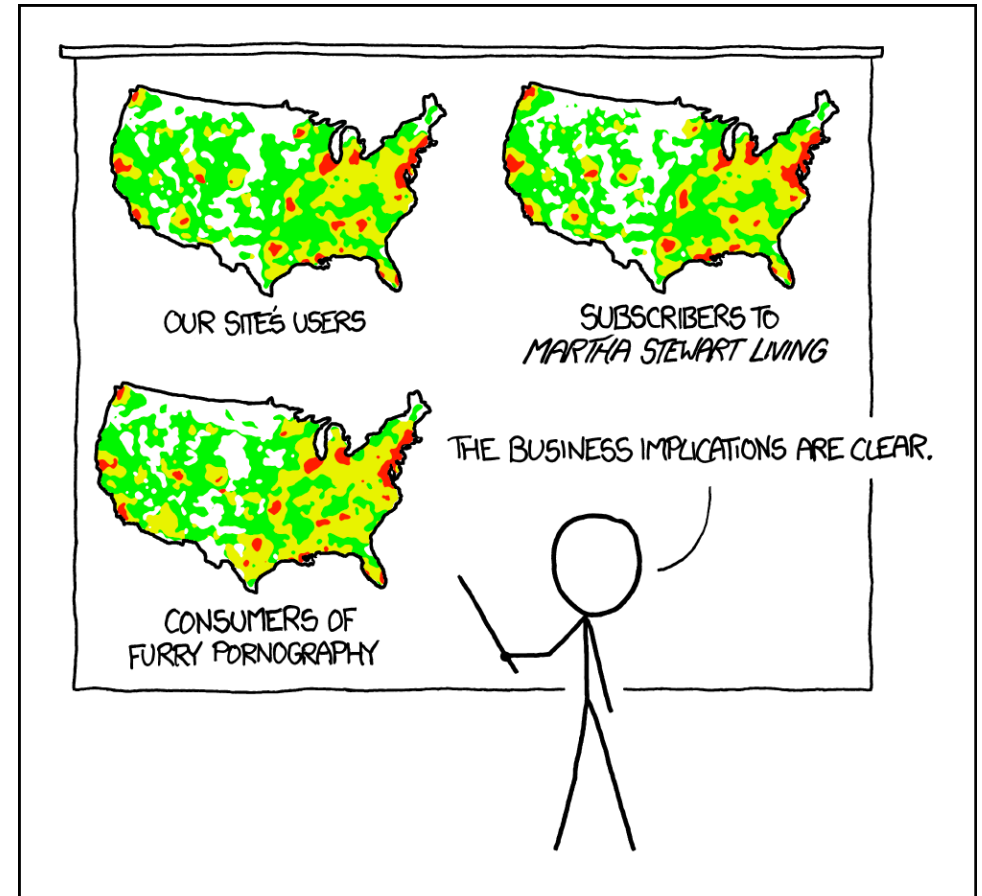


Misleading uses of area



The Problem with Absolute Values

- Sometimes, data should be relative



PET PEEVE #208:
GEOGRAPHIC PROFILE MAPS WHICH ARE
BASICALLY JUST POPULATION MAPS

Data/Ink Ratio

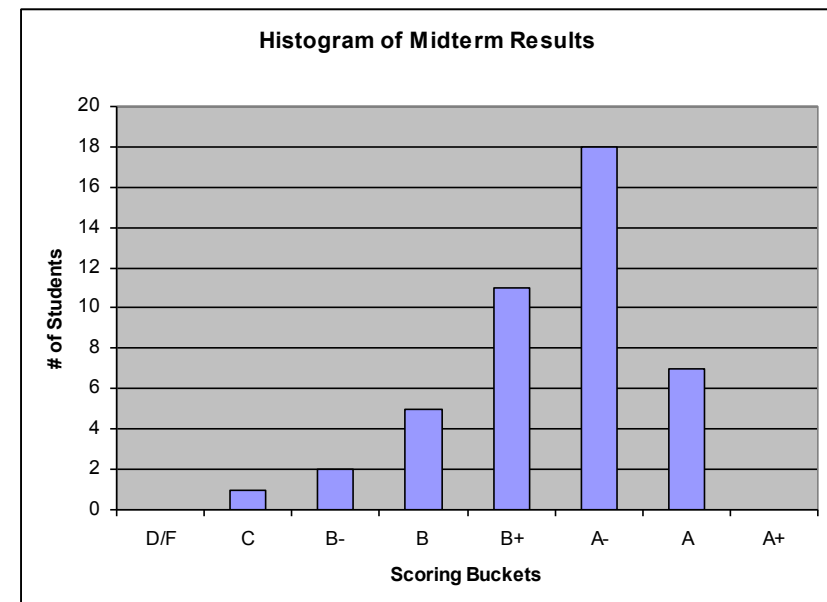
$$\text{Data-ink ratio} = \frac{\text{data-ink}}{\text{Total ink used to print graphic}}$$

= Proportion of a graphic's ink devoted to the non-redundant display of data-information.

< 0.001

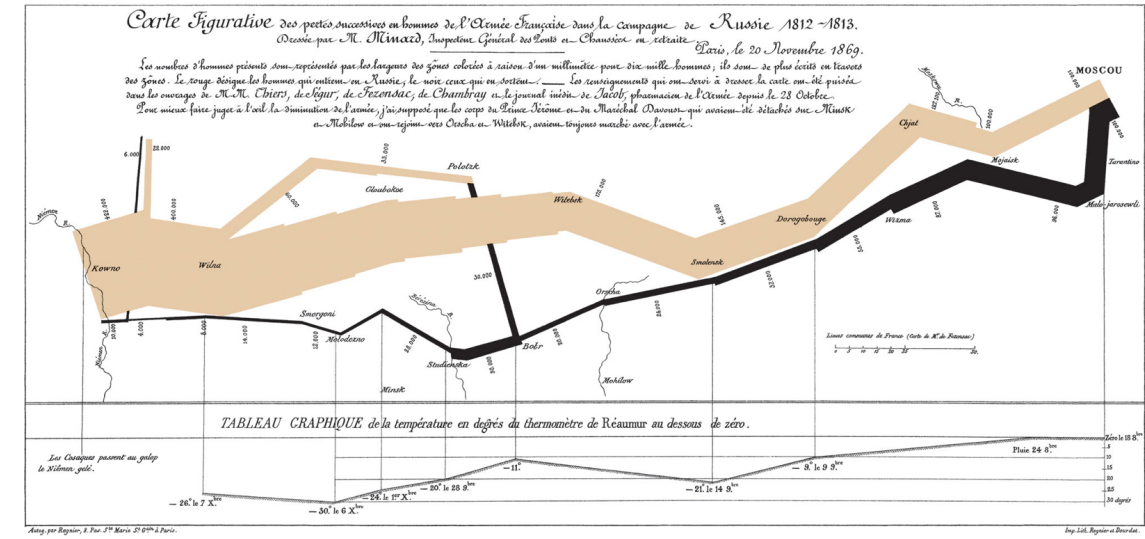


< 0.05 !!!



Helping Comparison

- Make the scale clear
 - Use intuitive scales
 - Make the baseline clear
- Provide context:
 - Examples
 - Generalization and Comparisons
 - Annotations and explanations



Carte Figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812-1813.

Dressée par M. Minard, Inspecteur Général des Ponts et Chaussées en retraite Paris, le 20 Novembre 1869.

Les nombres d'hommes présents sont représentés par les largeurs des zones colorées à raison d'un millimètre pour dix mille hommes; ils sont de plus écrits en travers des zones. Le rouge désigne les hommes qui entrent en Russie, le noir ceux qui en sortent. — Les renseignements qui ont servi à dresser la carte ont été puisés dans les ouvrages de M.M. Chiers, de Ségur, de Fezensac, de Chambray et le journal inédit de Jacob, pharmacien de l'Armée depuis le 28 Octobre.

Pour mieux faire juger à l'œil la diminution de l'armée, j'ai supposé que les corps du Prince Jérôme et du Maréchal Davout qui avaient été détachés sur Minsk et Mohilow et ont rejoint vers Orscha et Witebsk, avaient toujours marché avec l'armée.

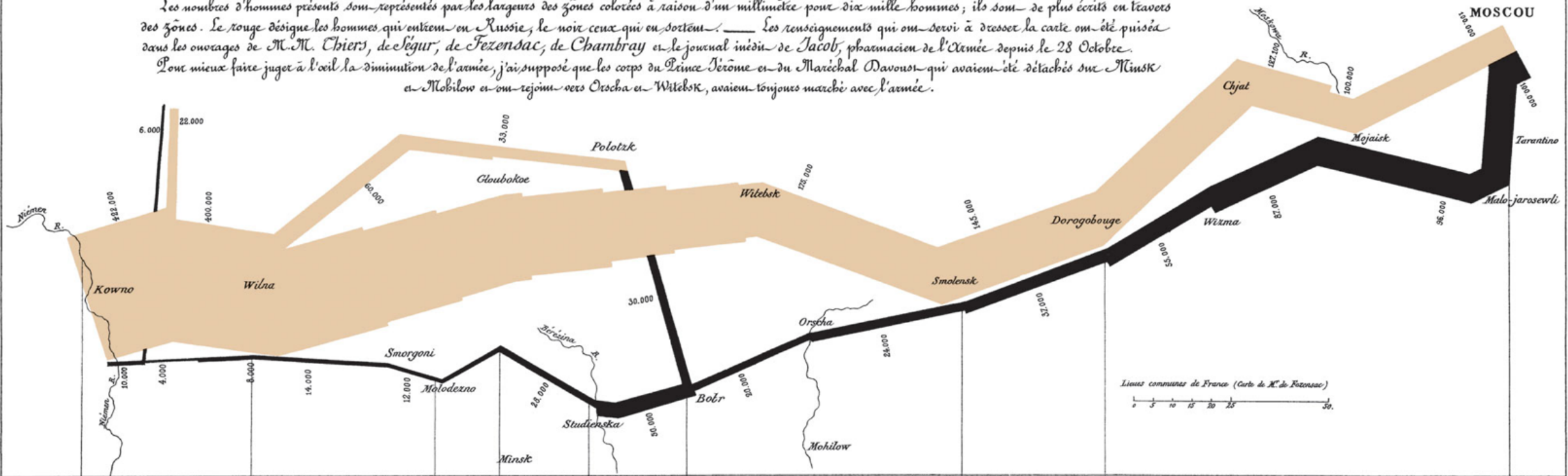
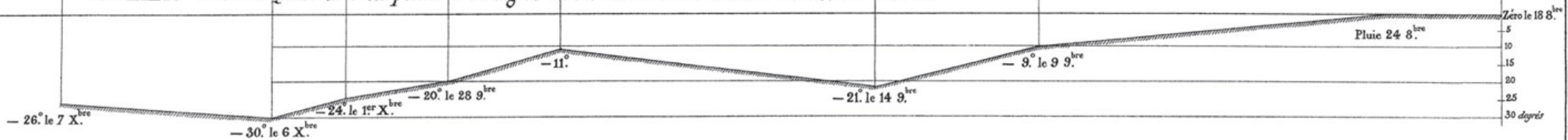


TABLEAU GRAPHIQUE de la température en degrés du thermomètre de Réaumur au dessous de zéro.

Les Cosaques passent au galop le Niémen gelé.

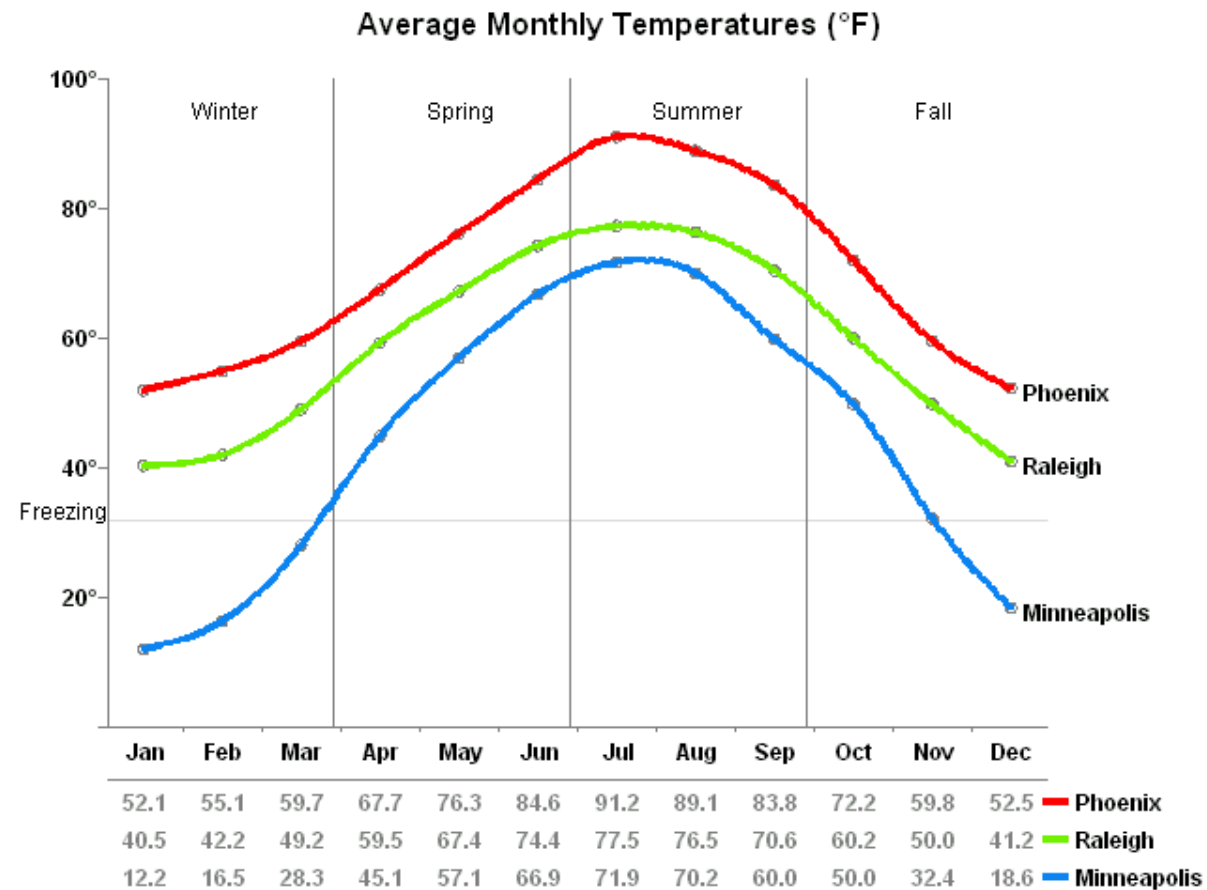


Auég. par Regnier, 8. Par. 5^{me} Marie 5^{me} G^{de} à Paris.

Imp. Lith. Regnier et Douvret.

The French engineer, Charles Minard (1781-1870), illustrated the disastrous result of Napoleon's failed Russian campaign of 1812. The graph shows the size of the army by the width of the band across the map of the campaign on its outward and return legs, with temperature on the retreat shown on the line graph at the bottom.

Make comparisons explicit



A great visualization can help you understand why the data is a certain

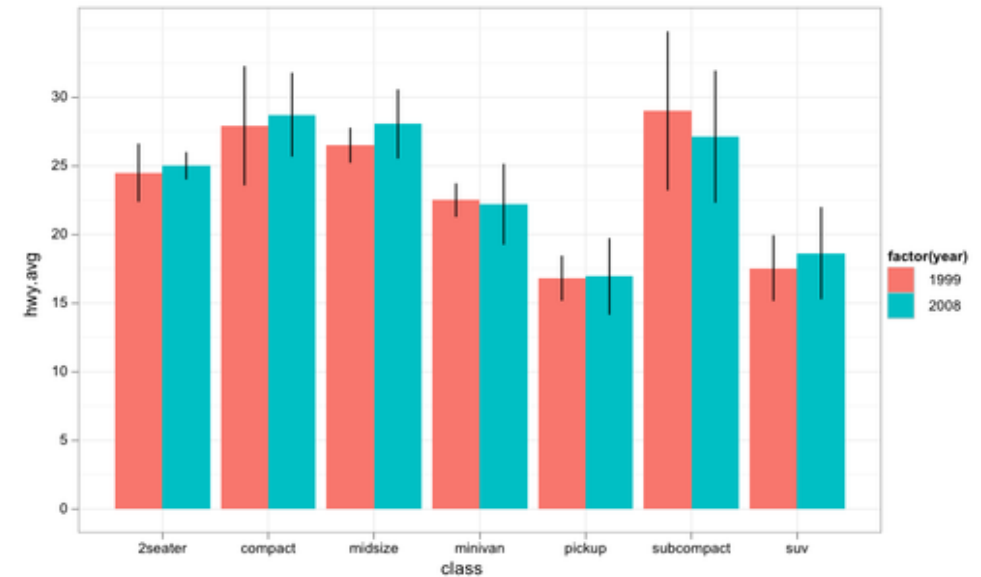
Baseline Example

Show baseline: The Dow Jones Average provides a baseline for comparison.

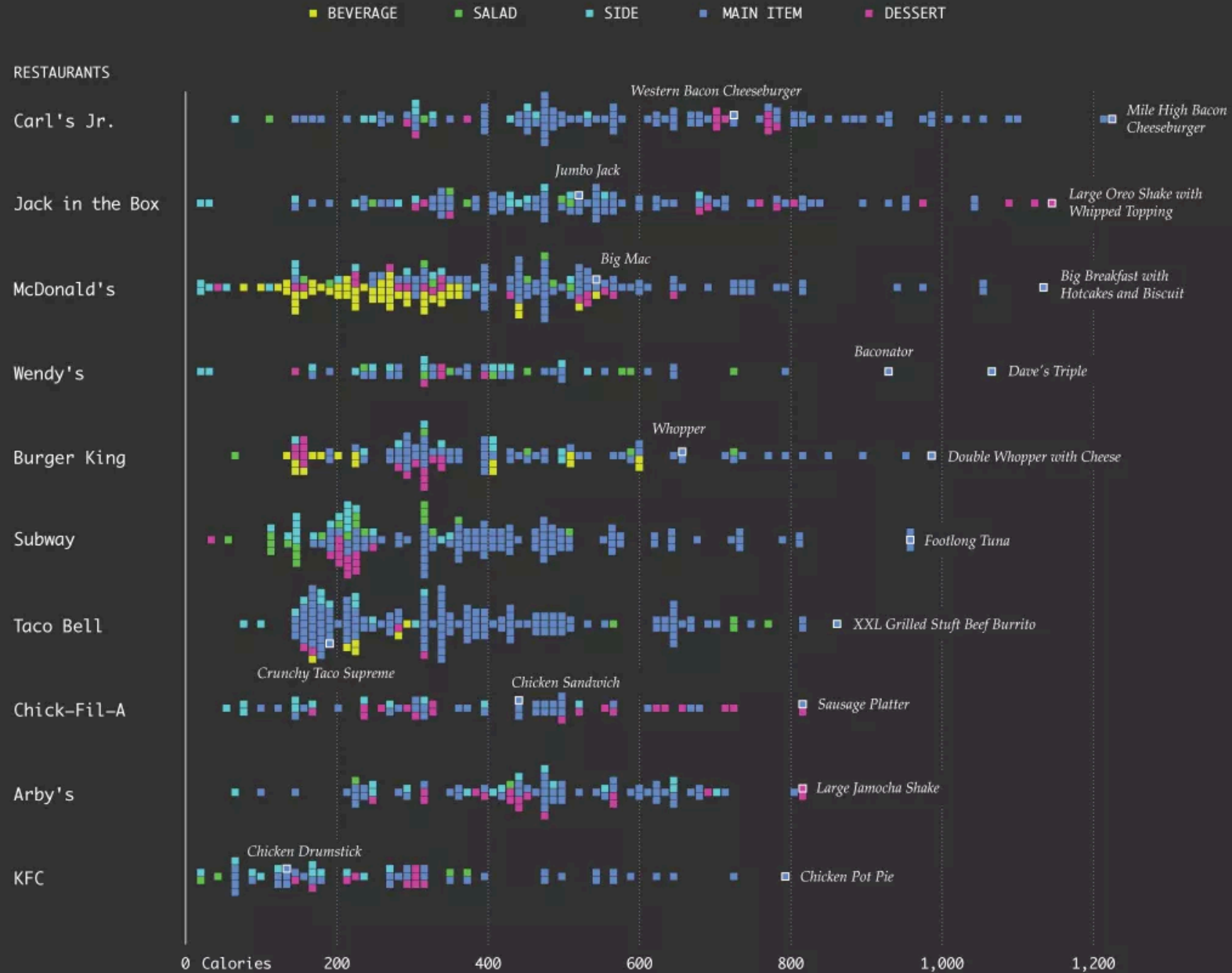


Providing Context: Distribution

- If bars are used to represent the average, error bars are necessary!
- The error bar visualize some form of variance
- Mostly standard error or standard deviation.

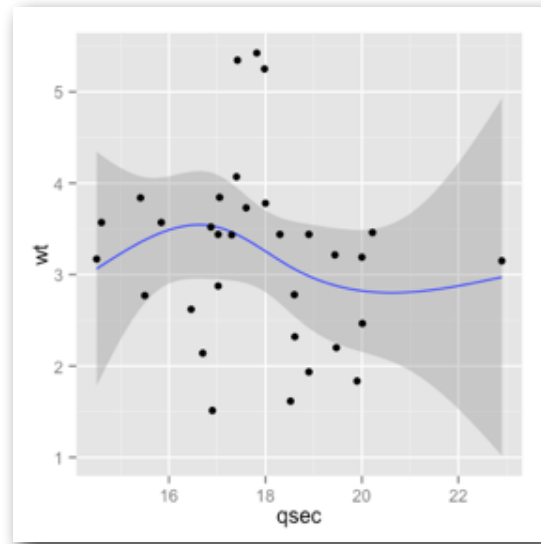


CALORIES IN FAST FOOD MENU ITEMS

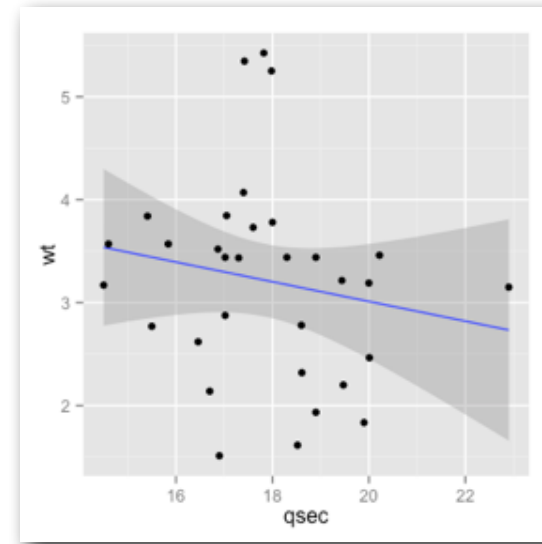


Context: Generalization

Adding generalizations to the graphs to provide aggregative data.



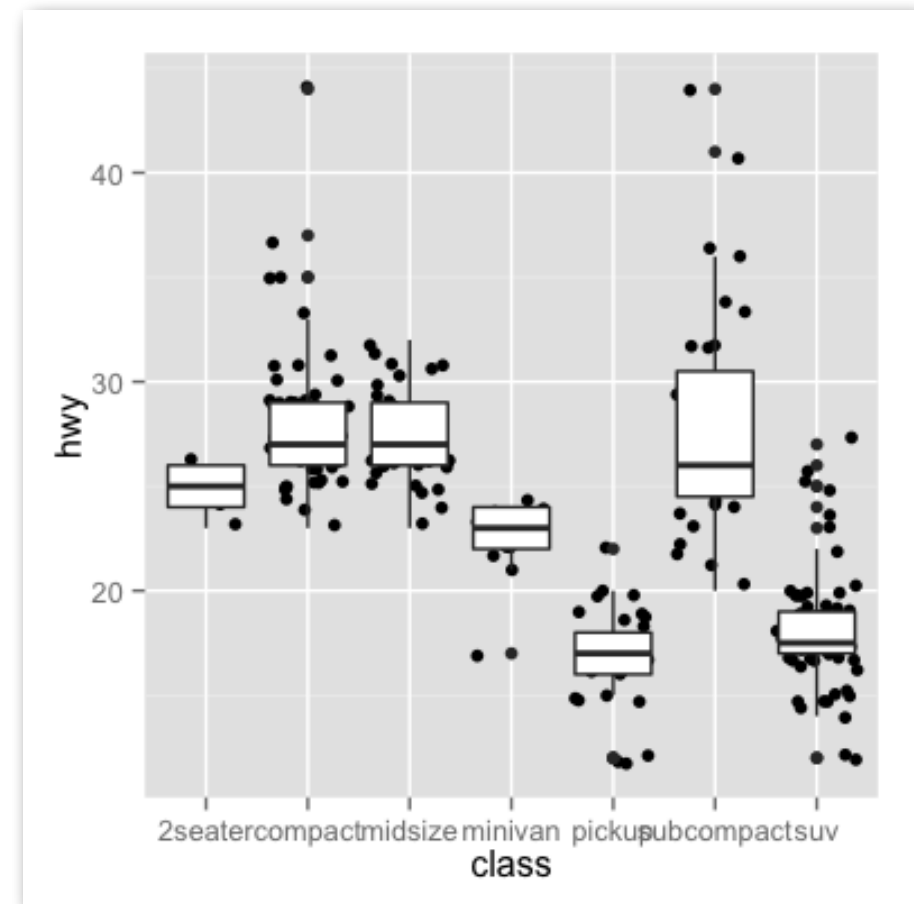
Local Regression Model



Linear Model

Context: Examples

Use example data to provide context to generalizations.



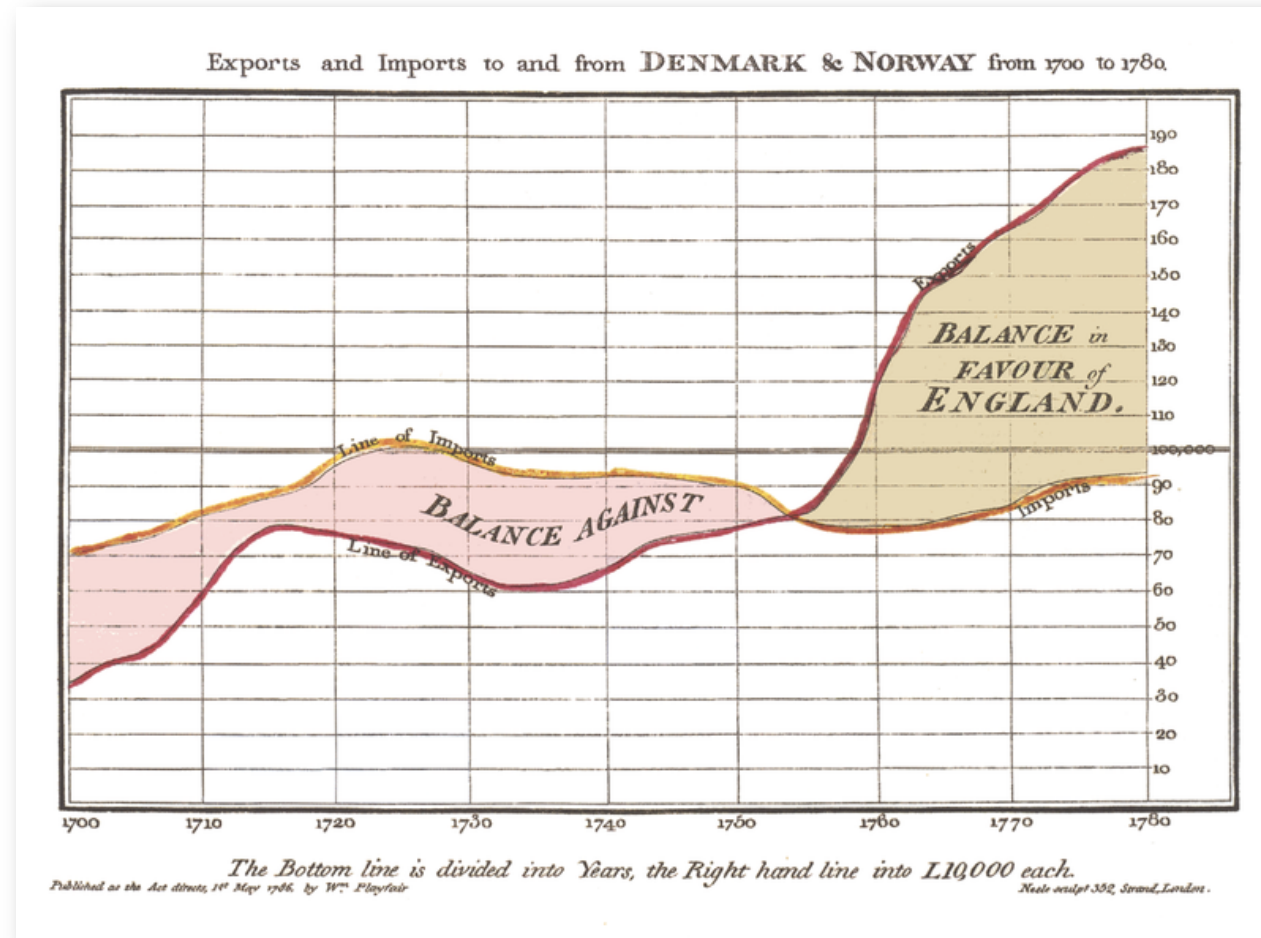
Annotations and explanations

Annotations allow people to understand the data by themselves, and draw conclusions from the data.

A1. Since October 2012 Apple stock has fallen on hard times, with increasing competition from Google, and a more fragmented market for smart-phones.



Time Series Visualization



Summary