# 01 Introduction

## 36-721 Statistical Graphics and Visualization
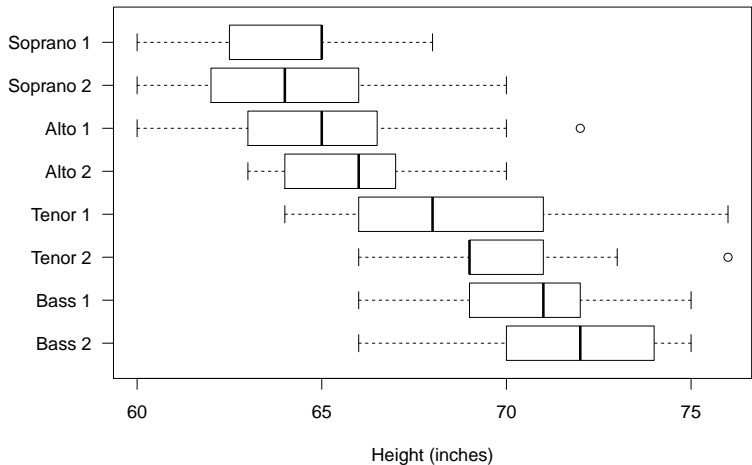
Jerzy Wieczorek

9/1/15

# Examples and context

What good is data visualization?

What can we aspire to?

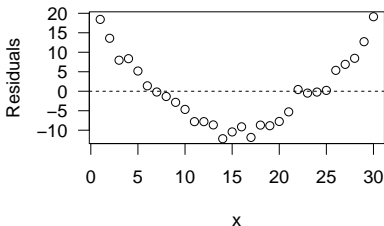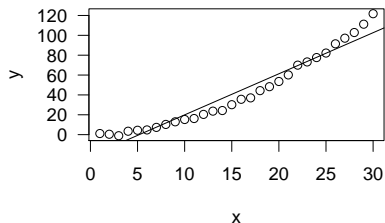How do statistical graphics fit in with other flavors of visualization and information design?
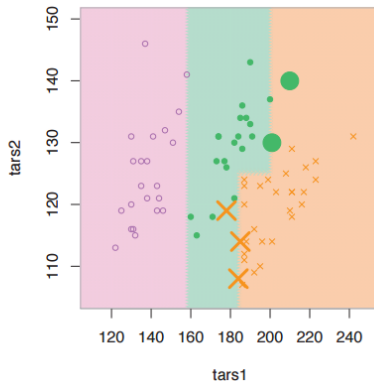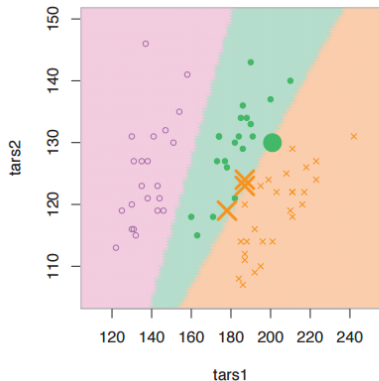
# Statistical graphics

## EDA

# Statistical graphics

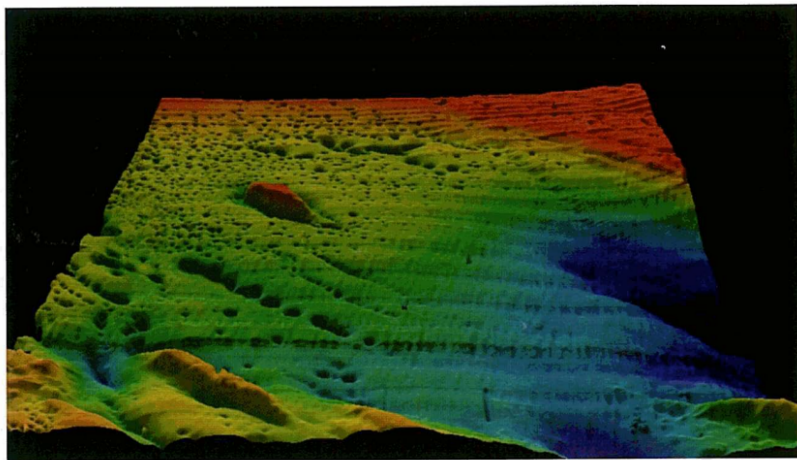Regression diagnostics

# Statistical graphics

**Classifier diagnostics**
Cook & Swayne, *Interactive and Dynamic Graphics for Data Analysis, With R and Ggobi*

# Scientific visualization
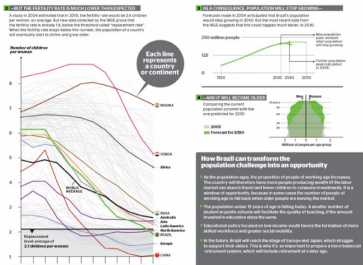
Ware, *Information Visualization*



Passamoquoddy Bay visualization. *Data courtesy of the Canadian Hydrographic Service.*

# Statistical graphics and scientific visualization

- See a huge dataset all at once
- Find interesting features at different scales
- Notice data anomalies
- Propose scientific hypotheses

# Infographics & data journalism

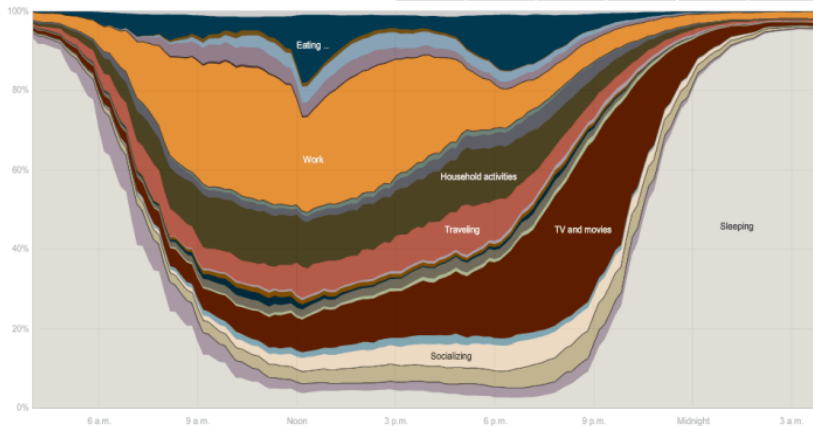Static infographic by **Alberto Cairo**

# Infographics & data journalism

Interactive graphic by **New York Times**

# Animation and narration

**Hans Rosling's TED talks**
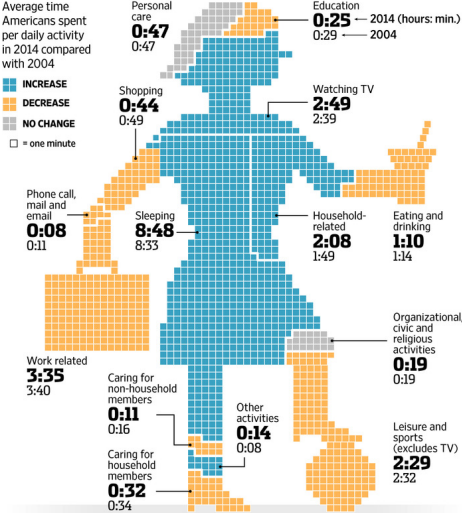
Watch **2:15-5:01 here**

# Data art

## Wall Street Journal



### A Day in the Life

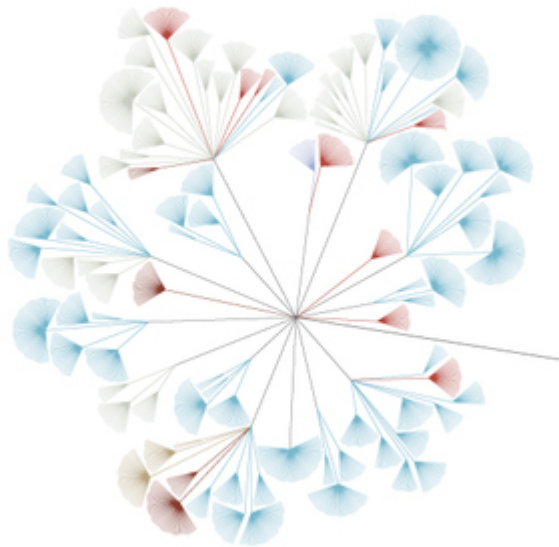Average time Americans spent per daily activity in 2014 compared with 2004

- **INCREASE** (blue)
- **DECREASE** (orange)
- **NO CHANGE** (gray)

□ = one minute

Personal care
**0:47**
0:47

Education
**0:25** ← 2014 (hours: min.)
0:29 ← 2004

Shopping
**0:44**
0:49

Watching TV
**2:49**
2:39

Phone call, mail and email
**0:08**
0:11

Sleeping
**8:48**
8:33

Household-related
**2:08**
1:49

Eating and drinking
**1:10**
1:14

Work related
**3:35**
3:40

Caring for non-household members
**0:11**
0:16

Other activities
**0:14**
0:08

Organizational, civic and religious activities
**0:19**
0:19

Leisure and sports (excludes TV)
**2:29**
2:32

Caring for household members
**0:32**
0:34

Note: Time may not total 24 hours due to rounding.
Source: Labor Department

Christopher Kaeser/THE WALL STREET JOURNAL.
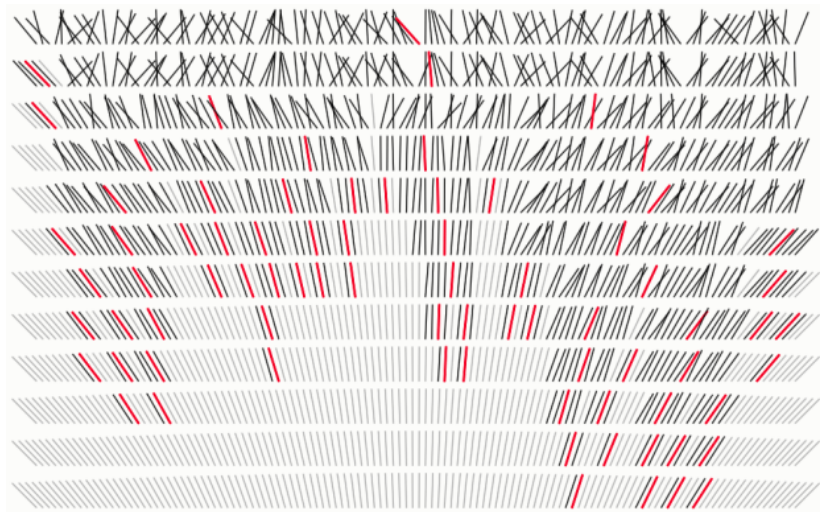
# Data art

## Stephanie Posavec

# Information design, visual explanation

**Wayfinding**

# Information design, visual explanation

**Instruction**

# Information design, visual explanation

**Visualizing algorithms**

# Course info

How is the class organized? How is it graded?

# How would you organize this course?

What topics to cover? How to group them?

# How I've organized this course

Course objectives (see Syllabus) are based on frameworks and principles

In class we will

- ▶ Look at graphics
- ▶ Discuss whether they work
- ▶ Propose reasons why they (don't) work
- ▶ Compare to frameworks/principles in the literature
- ▶ Critique graphs and make new ones using our principles
- ▶ Learn software tools as needed for implementation

# Assessments

Each assessment (homework, critique, project) targets one of the learning objectives. No points—just a rubric for each assignment.

Final grade depends on which assignments you completed and how well (see Syllabus).

You can revise and resubmit. But first submissions must be on time and show sincere effort! *(to keep grading manageable for us)*

# Syllabus

- Office hours and due dates
- Objectives
- Texts and software
- Assessments
- Administrivia
- Schedule

Be sure to note:

- Passing grade for CMU graduate students is B- or above.
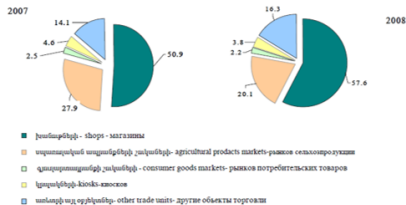- R software is required for Statistics MSPs.
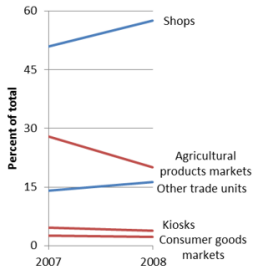
# Schedule of topics

What will we study?

# Readability and best practices

**Armenian National Statistical Service** and my own remake



TOTAL VOLUME OF RETAIL TRADE TURNOVER BY FORMATION SOURCES
ОБЩИЙ ОБЪЕМ РОЗНИЧНОГО ТОВАРООБОРОТА ПО ИСТОЧНИКАМ ФОРМИРОВАНИЯ



Total Volume of Retail Trade Turnover by Formation Sources
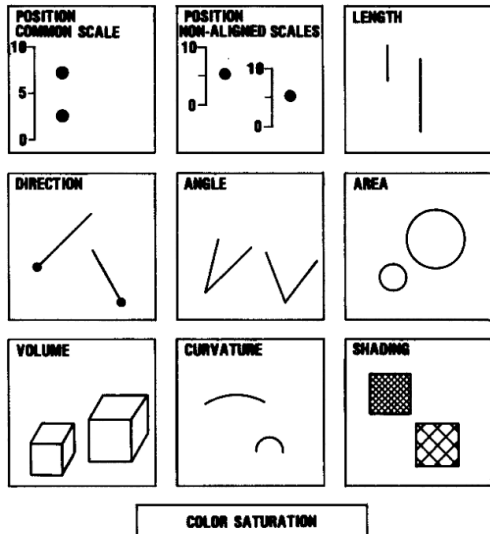
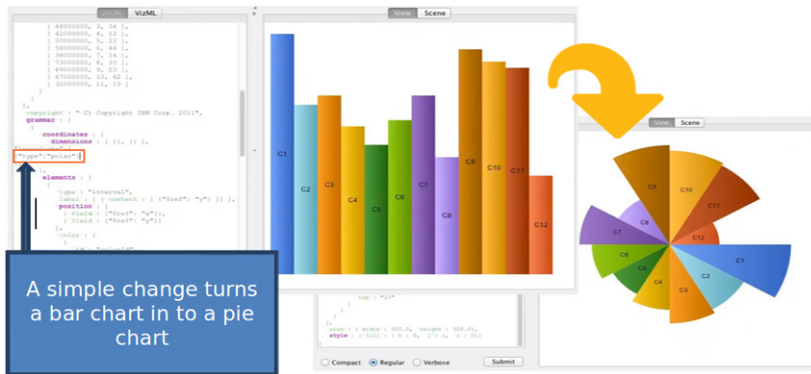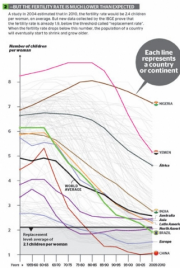# Human visual perception
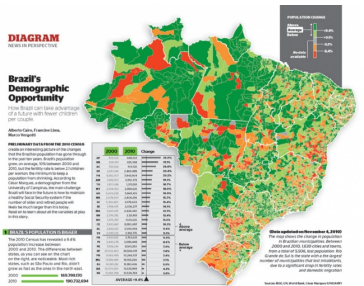
## Cleveland and McGill (1984)



Figure 1. Elementary perceptual tasks.

# The Grammar of Graphics

**IBM's VizJSON**, R's ggplot2, SPSS's GPL and Visualization Designer, Tableau…



A simple change turns a bar chart in to a pie chart
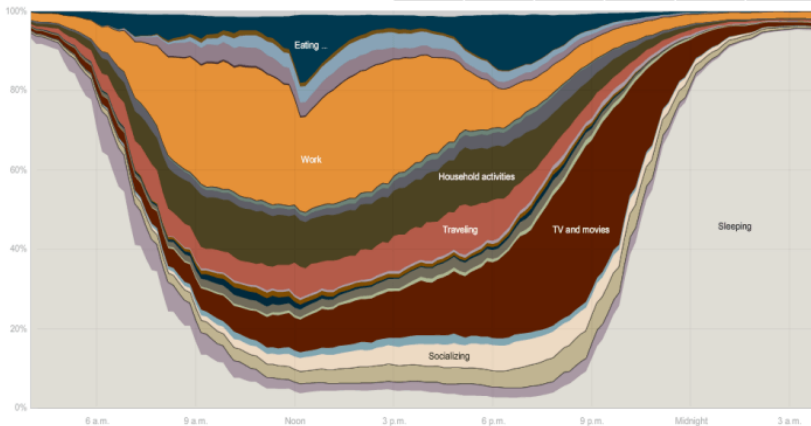
# Graphic design

# Interaction design



## Everyone
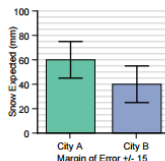Sleeping, eating, working and watching television take up about two-thirds of the average day.

| Everyone | Employed | White | Age 15-24 | H.S. grads | No children |
| Men | Unemployed | Black | Age 25-64 | Bachelor's | One child |
| Women | Not in lab... | Hispanic | Age 65+ | Advanced | Two+ children |

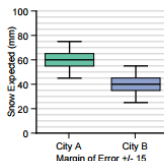# The dataviz research literature

**Correll and Gleicher (2014)**



Error Bars Considered Harmful:
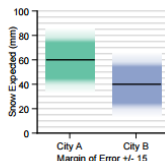Exploring Alternate Encodings for Mean and Error

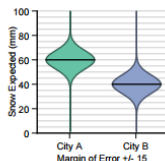Michael Correll *Student Member, IEEE*, and Michael Gleicher *Member, IEEE*

(a) **Bar chart** with error bars: the height of the bars encodes the sample mean, and the whiskers encode a 95% t-confidence interval.

(b) **Modified box plot**: The whiskers are the 95% t-confidence interval, the box is a 50% t-confidence interval.

(c) **Gradient plot**: the transparency of the colored region corresponds to the cumulative density function of a t-distribution.
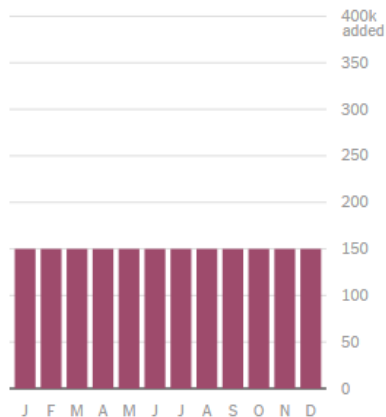
(d) **Violin plot**: the width of the colored region corresponds to the probability density function of a t-distribution.
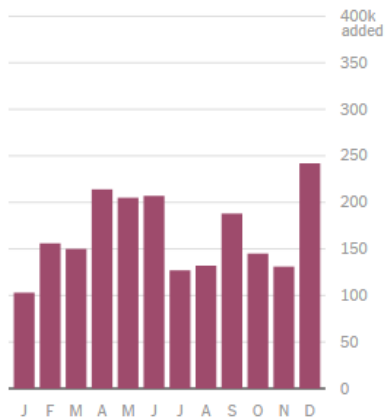
# Communicating statistical ideas

Sampling variation, via **New York Times**

# Graphics for statistical analysis

# Maps and cartography

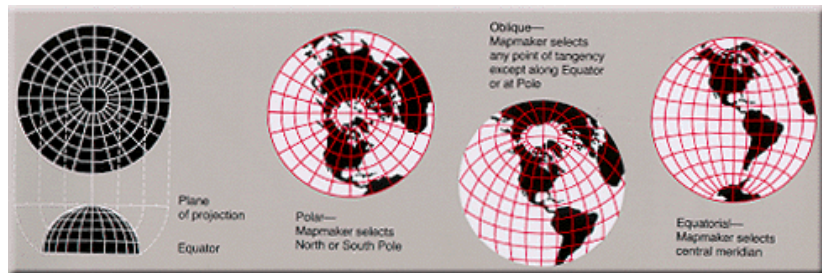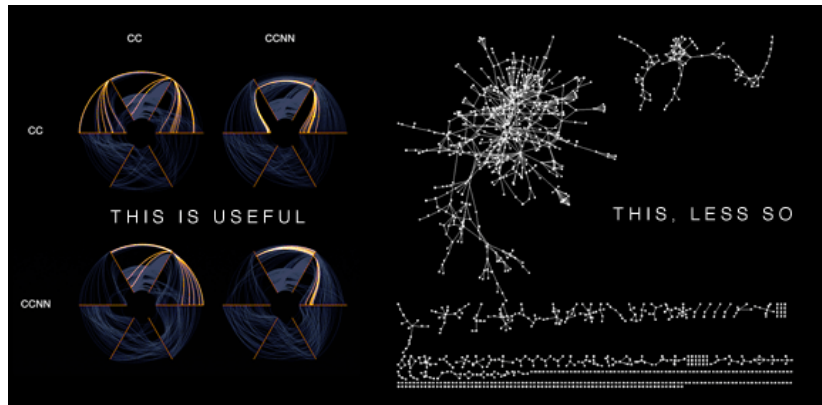**Map projections**, choropleths, cartograms

# Trees and networks

**Hiveplot network diagrams**

# TBD—ideas

- Vector fields
- Data sonification
- D3.js practice
- Chart zoo
- Table design
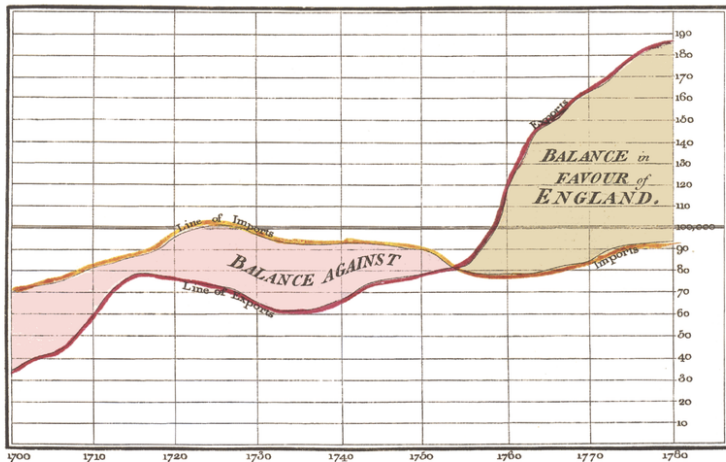- History of visualization
- *Your topic here?*

# Historical classics

So you can nod and say, "Oh yeah, I know that one"

# Playfair

**Time series**, 1786



Exports and Imports to and from DENMARK & NORWAY from 1700 to 1780

# Snow

**Cholera map**, 1854

# Nightingale

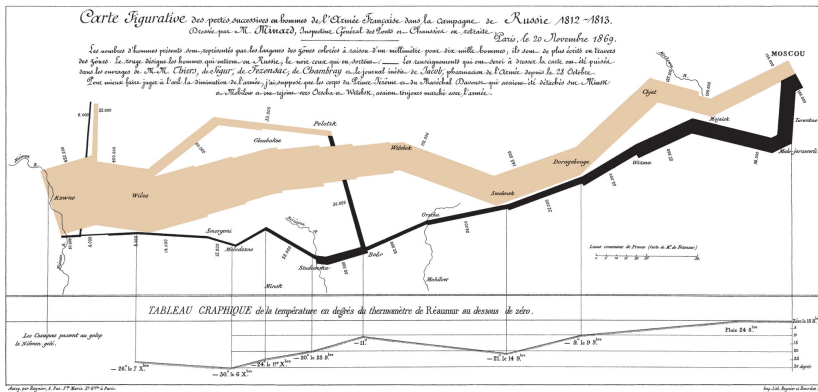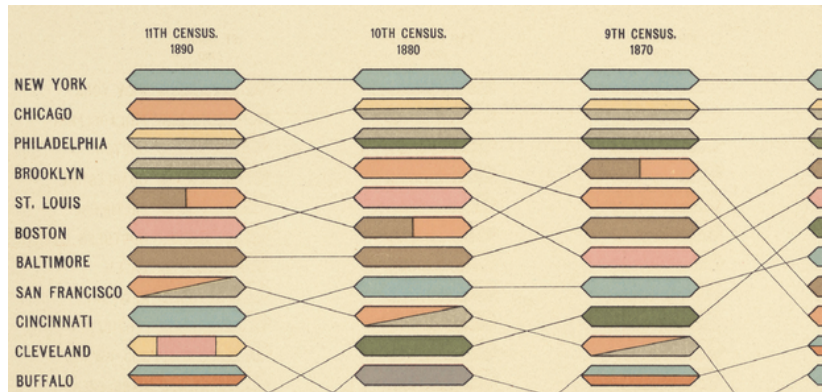**Polar area diagram**, 1858

# Minard

**Napoleon's march**, 1869

# Statistical Atlas of the United States

**Statistical Atlases** by the U.S. Census Bureau, 1870-1890, including **ranks of city populations 1790-1890**

# Neurath

**Isotype**, 1920s

# Anscombe

**Anscombe's quartet**, 1973 (via **Tufte**, *The Visual Display of Quantitative Information*)

Graphics *reveal* data. Indeed graphics can be more precise and revealing than conventional statistical computations. Consider Anscombe's quartet: all four of these data sets are described by exactly the same linear model (at least until the residuals are examined).
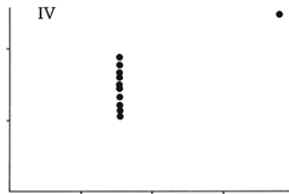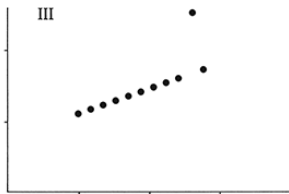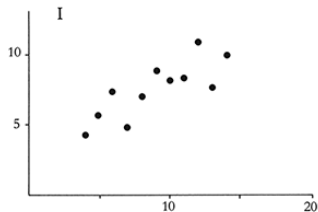
| I | | II | | III | | IV | |
|---|---|---|---|---|---|---|---|
| X | Y | X | Y | X | Y | X | Y |
| 10.0 | 8.04 | 10.0 | 9.14 | 10.0 | 7.46 | 8.0 | 6.58 |
| 8.0 | 6.95 | 8.0 | 8.14 | 8.0 | 6.77 | 8.0 | 5.76 |
| 13.0 | 7.58 | 13.0 | 8.74 | 13.0 | 12.74 | 8.0 | 7.71 |
| 9.0 | 8.81 | 9.0 | 8.77 | 9.0 | 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 |

$N = 11$
mean of X's = 9.0
mean of Y's = 7.5
equation of regression line: $Y = 3 + 0.5X$
standard error of estimate of slope = 0.118
$t = 4.24$
sum of squares $X - \bar{X} = 110.0$
regression sum of squares = 27.50
residual sum of squares of Y = 13.75
correlation coefficient = .82
$r^2 = .67$

# Anscombe

And yet how they differ, as the graphical display of the data makes vividly clear:
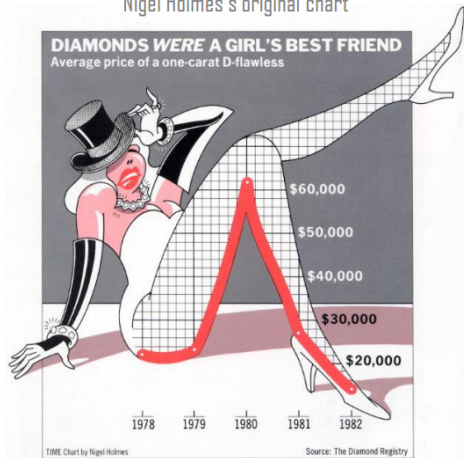
F. J. Anscombe, "Graphs in Statistical Analysis," *American Statistician*, 27 (February 1973), 17–21.
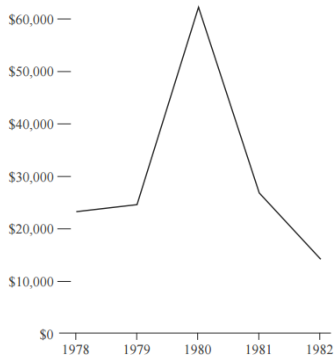
# Holmes vs. Tufte

Holmes, *TIME*, 1982(?); Tufte, *Envisioning Information*, p. 34;
Cairo, *The Functional Art*, p. 61-70
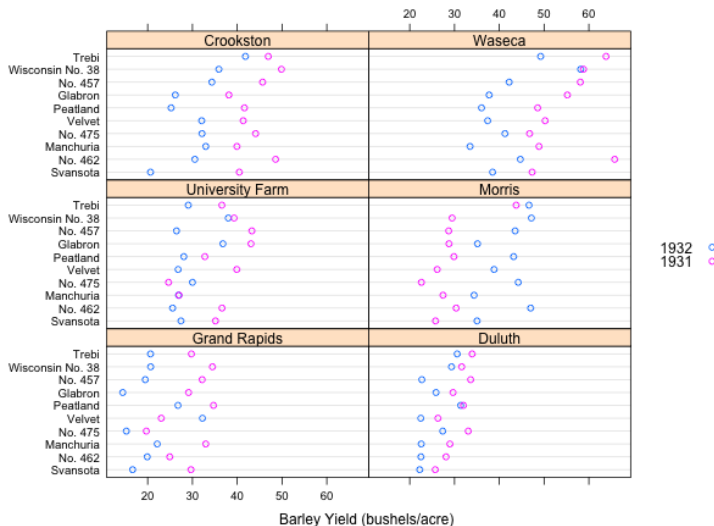
# Cleveland

Trellis dotplot of **barley yield data** (Cleveland, 1993), though **see rejoinder**

# Classic books

- Bertin, *Semiology of Graphics*
- Cleveland, *The Elements of Graphing Data* and *Visualizing Data*
- Wilkinson, *The Grammar of Graphics*
- Tufte, *The Visual Display of Quantitative Information* and *Envisioning Information*
- Wainer, *Visual Revelations* and others

# Next time

What to prepare? What'll we cover?

# Prepare for next time

- Install and test-drive your statistical graphics software
- Look at HW1 (data + rubric), let me know if unclear
- Readings: Cairo Ch 1-4; Donahue p. 1-23
- Blogs to follow
  - **Nathan Yau**
  - **Alberto Cairo**
  - **Robert Kosara**
  - **Kaiser Fung**
  - **Di Cook**

# Next time we'll cover

- Best practices for most common 1D/2D charts and tables
- Image formats, resolution, saving plots
- A few handy tricks (logs, loess, jitter)
- R users: bring laptops to follow along

# Software installation

Let's get everything installed and debugged to prepare for future classes.

# Software installation

- Who'll use R? What else will be used?
- R users: install and test-drive
    - R and RStudio
    - `ggplot2`, `knitr`, `shiny` packages
- Tableau users: install student license
    - **www.tableau.com/academic/students**
- Also consider
    - D3.js
    - Inkscape