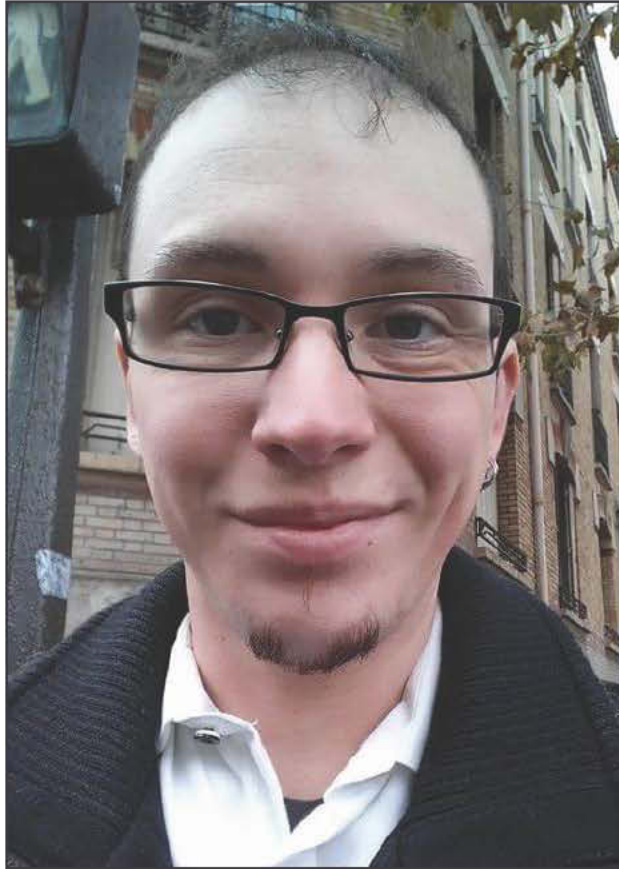MassMutual DSDP 2017:

# INTRODUCTION TO DATA VISUALIZATION

June 8, 2017

R. Jordan Crouser & Amelia McNamara

Statistical & Data Sciences

Smith College
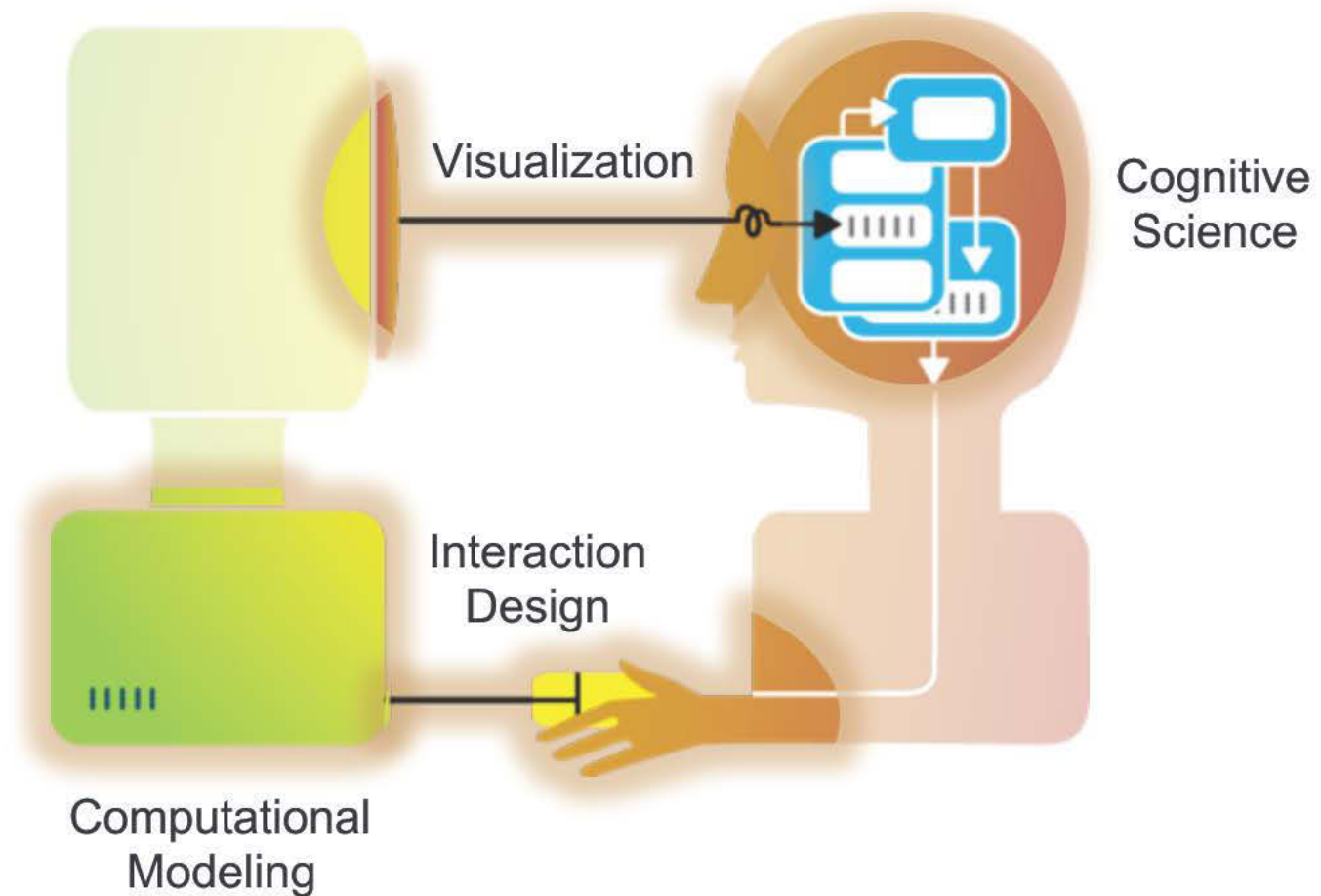
# People



### Jordan
(computer scientist)



### Amelia
(statistician)

# Our research (broadly)

# Housekeeping



jcrouser.github.io/datavis

# About this course



Visualization

# What is visualization?

# What is visualization?

# Perhaps a more helpful question:

What are some ways
a "visualization" can be **useful**?

# Does it help you spot trends?

# Does it help you explore?

# Does it tell a story?

# Visualization (def.)

Visual **representations** of data that reinforce human **cognition**

# Wait… what is "data"?

# Data: a definition

Data is a set of *variables* that capture various aspects of the world:



*Tuition rates, enrollment numbers,*
*public vs. private, etc.*

# Data: a definition

A dataset also contains a set of *observations* (also called *records*) over these variables. For example:



*tuition* = $46,288, *enrollment* = 2,563, *private, etc.*

# Data: a definition

A dataset also contains a set of *observations* (also called *records)* over these variables. For example:



*tuition* = $16,115,  *enrollment* = 28,635,
*public, etc.*

# One way to think about this:

# Another way to think about this

```
class school_obs:
    def __init__(tuition, enrollment,
                 pub_or_priv):
        self.tuition = tuition
        self.enrollment = enrollment
        self.pub_or_priv = pub_or_priv
```

VARIABLES

OBSERVATIONS

# Why is this important?

- Data have dimensions

- Visualizations have dimensions, too

- To build visualizations, we need to **map** data dimensions to visual dimensions

# Key question for this course

Which **data dimension** should be mapped

to which **visual dimension?**

# Answer: it depends



Average Height for Youth Sports Participants

# A quick history lesson…

# (Incomplete) History of Visualization: 15,000BC



15,000 BC.  Laxcaux, France

# (Incomplete) History of Visualization: 900s



*"De cursu per zodiacum", illustrator unknown*

# (Incomplete) History of Visualization: 1970s



- CAD/CAM, building cars, planes, chips
- Starting to think about: 3D, animation, edu, medicine

# (Incomplete) History of Visualization: 1980s



- Scientific visualization, physical phenomena
- Starting to think about: photorealism, entertainment

# (Incomplete) History of Visualization: 1990s



- Information visualization, storytelling
- Starting to think about: online spaces, interaction

# (Incomplete) History of Visualization: 2000s



- Coordination across multiple views, interaction
- Starting to think about: sensemaking, provenance

# Discussion: what are they all trying to do?

# Visualization helps shape *mental models*

# Information overload

- We are exposed to huge amounts of information all the time

- So much, in fact, that we can't process it all fast enough

# Mental models

To cope, we construct **mental models:** abstracted, simplified versions of the world that are more manageable

# Mental Models: a Sketch

# 1. We tend to see what we expect to see

# 2. Mental models form quickly, & update slowly

# 3. New information gets incorporated into the existing model

# 4. Initial exposure interferes with accurate perception



**Blur size**

128px
64px
32px
16px
8px
None

# The good, the bad, and the ugly…

The good:
- Well-tuned mental models let us process information quickly
- Frees up more processing power to synthesize information

The bad:
- People (esp. experts) tend not to notice information that contradicts their mental model
- A "fresh pair of eyes" can be beneficial

The ugly:
- Mental models are unavoidable: everyone has them, and they're all different
- **Key:** be aware of how mental models form, how they shape perception, and how to support (or challenge) them

# So what do we have to work with?

# Graphical primitives

The images we draw are composed of marks: like ink

# Visual dimension: position

- Encode information using **where** the mark is drawn
- Some examples:

# Visual dimension: size

- Encode information using **how big** the mark is drawn
- Examples:

# Visual dimension: value

- Encode information using **how dark** the mark is drawn
- Example:



Legend:
- 0 - 1
- 1 - 3
- 3 - 6
- 6 - 10
- 10 - 16
- 16 - 30
- 30 - 85
- 85 - 160
- 160 - 550
- 550 - 1,100
- 1,100 - 2,500
- 2,500 - 5,000
- Over 5,000

# Visual dimension: color

- Encode information using the **hue** of the mark
- Examples:

# Visual dimension: orientation

- Encode information using how the mark is **rotated**
- Examples:

# Visual dimension: shape

- Encode information using how the mark is **shaped**
- Examples:

# Discussion

What makes a **good** encoding?

# Principle 1: expressiveness

- Encodes **all** the facts
- Example:

# Principle 1: expressiveness

- Encodes **only** the facts
- Example:



Adapted from Mackinlay J (1986) Automating the design of graphical presentations of relational information.

# Principle 2: consistency

- Use **consistent axes** when comparing charts



*misleading*          *improved*

Raina SZ, et al. (2005) Evolution of base-substitution gradients in primate mitochondrial genomes. Genome Res 15: 665-673.

M. Krzwinski, behind every great visualization is a design principle, 2012

# Principle 2: consistency

- A note on **legends**: order items according to appearance



consistent   inconsistent

| | consistent | | inconsistent |
| --- | --- | --- | --- |
| ☐ | A | ■ | A |
| ☐ | B | ◪ | B |
| ◪ | C | ☐ | C |
| ■ | D | ☐ | D |

# Principle 2: consistency

- Visual variation should **reflect and enhance** the underlying variation in the data
- Avoid **visually similar** encodings for independent variables
- Example:

# Principle 2: consistency

- Uniform size and alignment reduces visual complexity and aids interpretation
- Example:



*variation refactored*

Fig. 1: Sharov AA et al. (2005) Genome-wide assembly and analysis of alternative transcripts in mouse. Genome Res 15: 748-754.

Fig. 2: M. Krzwinski, behind every great visualization is a design principle, 2012

# Tufte, 1983

"Above all else, show the data."

The Visual Display
of Quantitative Information

EDWARD R. TUFTE

# Tufte, 1983

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

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

= 1 - proportion of a graphic that can be erased

# Tufte: maximize the data-ink ratio

# Familiar example

# Discussion

- What do you think of the data-ink ratio?
- Consider ways to **maximize** it…

# Principle 3: importance ordering

- Avoid unnecessary containment and repetition
- Example

# Principle 3: importance ordering

- Navigational aids shouldn't compete with data
- Avoid: heavy **axes**, **error bars** and **glyphs**



**lighest useable**

**darkest useable**

sparse

dense

15%

25%

45%

Heer J, Bostock M (2010) Crowdsourcing graphical perception: using mechanical turk to assess visualization design. Proceedings of the 28th international conference on Human factors in computing systems. Atlanta, Georgia, USA: ACM. pp. 203-212.

# Principle 3: importance ordering

- Simplify, simplify, simplify…



chart junk

visually concise

Sharov AA, et al (2006) Genome Res 16: 505-509.
Peterson J, et al. (2009) Genome Res 19: 2317-2323.
Thomson NR, et al. (2005) Genome Res 15: 629-640.
DB, Ko MS (2005) Genome Res 15: 748-754.

M. Krzwinski, behind every great visualization is a design
principle, 2012

# A caveat: "chart junk" and recall



**MONSTROUS COSTS**
**Total House and Senate campaign expenditures, in millions**

Bateman et al. "Useful Junk? The Effects of Visual Embellishment on Comprehension and Memorability of Charts", CHI 2010

# Chart junk and eye gaze



Figure 9. Percentage of on-screen time spent looking at different chart elements for Holmes and Plain charts.

# Lab 1: Deconstructing Data Graphics

- Break into groups of 2-3 people, and go to:

  [jcrouser.github.io/datavis/lab1.html](jcrouser.github.io/datavis/lab1.html)

- During this lab, we'll **critique** some professionally-made visualizations using these principles

- Try to think about the following questions:
  - What is the **first thing you notice** about this visualization?
  - What **point** is this visualization trying to make?
  - Who is the **intended audience**?
  - What is the visualization **doing well**?
  - What **problems** do you see with the visualization design?
  - **Why** do you think the designer made those choices?

# A Day in the Life

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

- **INCREASE**
- **DECREASE**
- **NO CHANGE**

☐ = 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

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

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

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

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

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

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

Christopher Kaeser/THE WALL STREET JOURNAL.

# What your BRAND COLOR

## SAYS ABOUT YOUR BUSINESS

**85%** Of Shoppers place color as a primary reason for why they buy product

**80%** of the brand recognition can be enhanced by color

People see color before they absorb anything else

**92%** believe that customers, color gives an image of impression to the quality

Color can improve readership by 40%, learning by 55-78% and 73% by comprehension

Upto 90% of assessments are based on color alone

**90%** believe that customers can get attracted with the color in choosing

**81%** think that color gives them a competitive edge

**76%** believe that the use of color makes their business appear larger to clients

**90%** believe that customers remember presentations and documents better when color is used.

**83%** believe that color makes them appear more successful

**84%** Say that color is the Primary reason to buy a particular product

2016 Iowa Democratic Presidential Caucus

| | |
|---|---|
| Hillary Clinton | 45.5% |
| Bernie Sanders | 43.1% |
| Martin O'Malley | 4.4% |

Sept. 30, 2015

| | |
|---|---|
| Clinton | 48.1% |
| Sanders | 31.5% |
| O'Malley | 3.6% |

HUFFPOST ➤ POLLSTER

July
2015

Jan.
2016

# EVENTS CONTRIBUTING TO DROP OF EURO

Investors think subprime would be the U.S. only crisis
**$1.4738**

Investors realize EU is also vulnerable to subprime
**$1.3919**

The decision of European Central Bank to increase the interest rate backfires
**$1.2545**

Investors are worried about the weakness of the EU economy
**$1.20**

USD per 1 EUR

Greek debt crisis flares up
**$1.2640**

Euro plummets as the Greece debt crisis worsens
**$1.2149**

The Ukraine crisis starts to heat up
**$1.36**

GREXIT

Fear that election may result in Greece leaving Eurozone
**$1.21**

ECB's purchase of eurobonds devalues the Euro
**$1.06**

Traders push Euro down and invest in Pound
**$1.0479**

2007   2008   2009   2010   2011   2012   2013   2014   2015

UNDER PRESIDENT OBAMA,
**MORE STUDENTS ARE EARNING** THEIR HIGH SCHOOL DIPLOMAS THAN EVER BEFORE

HIGH SCHOOL GRADUATION RATE

75% — 2007-08
75% — 2008-09
78% — 2009-10
79% — 2010-11
80% — 2011-12
81% — 2012-13
82% — 2013-14

#LeadOnEducation

SOURCE: U.S. DEPARTMENT OF EDUCATION, NATIONAL CENTER FOR EDUCATION STATISTICS

**57%**

of Europeans are worried their personal information is not safe.

Symantec.

# Illinois

Variable: Net Job Creation (Per 100)
Employees, Same Sex and Age Group
Year: 2000 Quarter:1
Sex: All and Age Group: Ages 19–21



Midpoint of Range: ▨ –5  ▩ –3  ▨ –1
                   ▨ 1   ▨ 3   ☐ 5

**Fig. 5.7** Job creation for young workers, by county, Illinois

# Who do Nike sponsor?

## International sports and events sponsor

The American based company is the largest sports supplier in the world, suppling equipment, shoes and apparel.
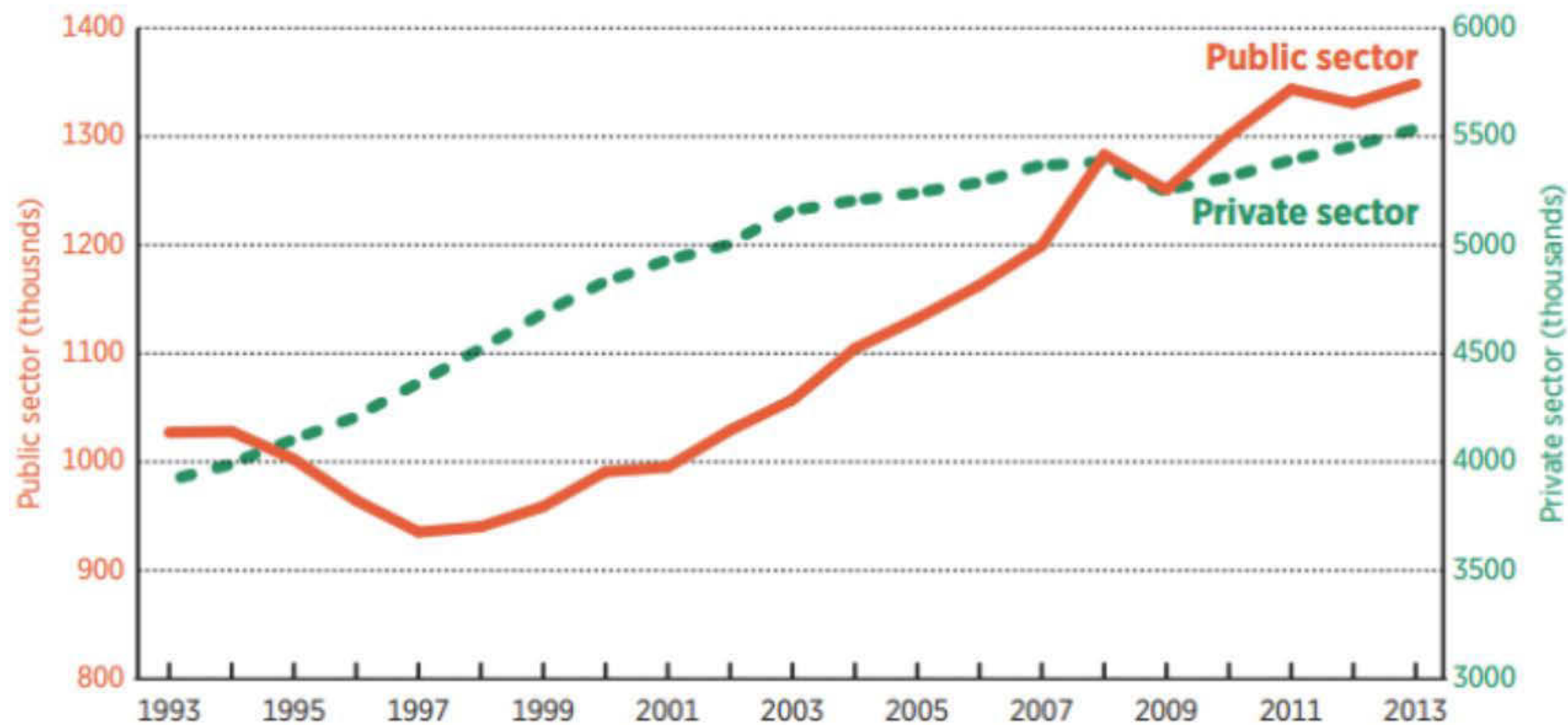
**1,016** athletes sponsored worldwide

**$230m** spent on athlete sponsors

**485** Male Footballers

**93** Basketball Players

**120** Male Track and Field Athletes

**64** Ice Hockey Players

**72** Female Track and Field Athletes

**29** American Footballers

**34** Baseball Players

**25** Male Golfers

**21** Cricketers

**31** Male Tennis Players

**1** Figure Skater

**10** Swimmers

**1** Martial Artist

**4** Female Golfers

**21** Female Tennis Players

**2** Cyclists

**2** Skiers

**15** Female Footballers

**2** Squash Players

**8** Rugby Players

**8** Boxers

# Figure 10: Public- and private-sector jobs (000s) in Ontario, 1993–2013

**2011**
**193,600**

**2015**
**117,161**

Despite the hysteria, the number of refugees in the UK has actually fallen by 76,439 since 2011.

Sample Description

Age
- 18-29
- 30-49
- 50-59
- 60+

Gender*
51 · 49

Educational level
- high
- medium
- low

10

* Due to rounding, numbers presented in this document may not add up precisely to 100%

## EPL (English Premier League)
- 6.3M
- 840K

## NBA
- 2.9M
- 620K

## NFL
- 2M
- 380K

## MLB
- 1.3M
- 320K

## NHL
- 600K
- 250K

# Coming up next

- Grammar of graphics
- Introduction to `ggplot2`
- Lab: Make a Scatterplot