# STAT 209
# Data Computing and Visualization

February 5, 2018

Colin Reimer Dawson

# Outline

"Data Science"

Intros

Some Terminology

Course Outline
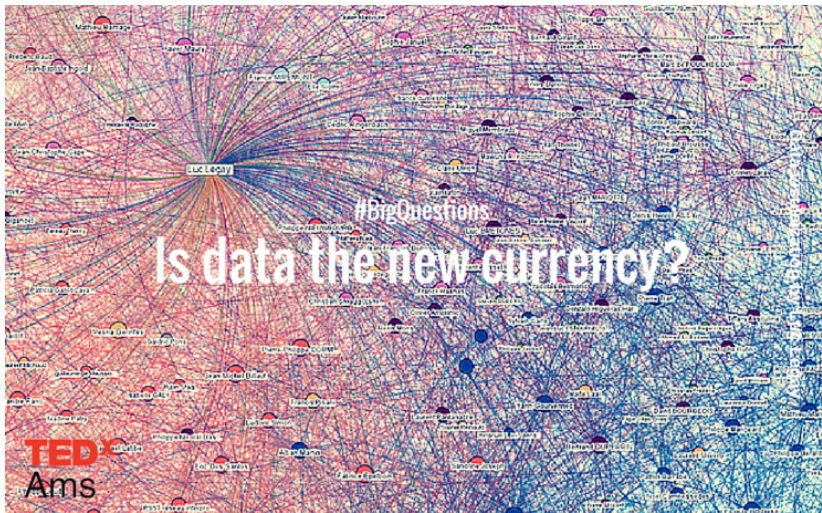
Data is the new black

"DATA IS THE NEW GOLD"

#BigQuestions
Is data the new currency?

# 'DATA' IS THE NEW CURRENCY

## WHAT OPPORTUNITIES ARE YOU MISSING?

# Some Cool Things you can do with data



Thanks to David Shuman at Macalester College for this slide

# DATA VISUALIZATION

RELIABLE, TIMELY, CONTENT

**DATA**

UNDERACHIEVING

INVISIBLE

POTENTIAL

THEME, COLOR, FONTS, READABLE

**DESIGN**

BORING

**GOOD INFOGRAPHIC**

AMATEUR

**STORY**

PROBLEM, CLEVER, MESSAGE, SOLUTION

DAMAGING

LIABILITY

EMBARRASSING

**SHAREABILITY**

VIRALITY, SEO, LOCATION, SOCIAL

© D. Zwiri

# Brainstorm

What is the difference between "data" and "information"?

# Cases

**Cases** When we collect data, we write down some measurements or characteristics of our **cases** — the individual "entities", sometimes called "observational units", that make up our dataset.

- The people in a survey or research study
- Plots of land in an agricultural experiment
- Days, in a weather dataset

# Categorical vs. Quantitative Variables

For each case we record one or more **variables**. One of the most basic distinctions is between **categorical** (or "qualitative") and **quantitative** data.

**Categorical:** "Qualitative" variable that divides cases into groups

**Quantitative:** Measures something on a scale; arithmetic makes sense

# Data Frames

A standard form for a dataset is a grid, called a **data frame**, where each row is a *case*, and each column is a *variable*.

| ID | Major | Height |
|----|-------------|--------|
| 1 | Neuroscience | 67 |
| 2 | CS | 71 |
| | ... | |
| 21 | Economics | 64 |

# Deconstructing Visualizations

For each of the following visualizations:

1. What are the cases (think "rows" of a dataset)?
2. What variables are depicted (think "columns" of a dataset)?
3. What graphical element (position, color, etc.) is used to encode each variable?

# Julia Louis-Dreyfuss



**Julia Louis-Dreyfus is good at almost everything**
IMDb ratings for appearances by Louis-Dreyfus

# Baseball Hits



**The sweet spot**
Scoring value (LWTS) of batted balls based on launch angle and speed off the bat, 2015 MLB

# Global Economic Growth



**Growth Across the Globe**

For the first time since the financial crisis a decade ago, all of the world's major economies are growing.

Canada +2.5%

U.K. +1.5%

Russia +1.8%

Euro Area +2.2%

U.S. +2.3%

Turkey +4.0%

China +6.6%

Japan +1.4%

Mexico +2.5%

India +6.2%

South Korea +2.7%

Indonesia +4.9%

Brazil +1.0%

Australia +1.8%

NO DATA  −7%  0  +1.8%  +3.6%  +5.6%  +8.4%

**Economic growth in 2017**

Year-over-year change in gross domestic product

China | United States | Euro area | India | Japan | Russia | Indonesia | Brazil | United Kingdom | Mexico | Turkey | South Korea

'00 '17

Some figures are estimates

Source: The Conference Board; Bureau of Labor Statistics | By Karl Russell

# Course Outline

- Part I: Basic Visualization (about 3 weeks)
- Part II: Data "Wrangling" (about 4 weeks)
- Part III: Dealing with "large" datasets (about 3 weeks)
- Part IV: Visualizing data with complex structure (spatial and text data) (about 2 weeks)

# On the web

- Course Website: `http://colindawson.net/stat209`
  - Syllabus, schedule, homework, slides, code, etc., there
- Blackboard: only for things that need a login/password protection
  - HW Solutions (when applicable)
  - Electronic submission of (some) assignments
- Slack: `stat209s2018.slack.com`, or download the app to computer/mobile device
  - Convenient one-stop place for all course-related electronic communication
- DataCamp
  - Many interactive tutorials to learn/practice computing tools
  - First couple of homework assignments there
- GitHub (later): Good way to track code changes/share code

# Graded Components

Course grade based on:

- ▸ Homework sets (15%)
- ▸ Labs (15%; completion/effort only)
- ▸ Four group visualization projects, one for each major "unit" (15% each)
- ▸ Participation/engagement (10%)
- ▸ Scheduled final exam day used for presentations of project 4

See the syllabus for Honor Code guidelines

# Structure of Class

- About half lecture/full group activities, half labs in small groups
  - Some labs are guided exercises to learn new tools
  - Final lab in each unit will ask you to "reverse engineer" a specific visualization from the web
- First three projects:
  - In class "workshop day", to work out the kinks together
  - In class "short presentation day", for feedback on a draft
  - Final writeup due a few days after that
- Project 4: longer presentation of a polished version during finals week

# This Week

- First homework is to complete a DataCamp chapter on using RStudio
- Wednesday: Lab 1 to get comfortable with R/RStudio
- Friday: Start on basic elements of visualization