# INTRO TO DATA SCIENCE
# LECTURE 1: DATA EXPLORATION

# WELCOME!

Course Website: https://www.schoology.com/

Instructors: Thomson Nguyen, Jacob Bollinger

E-mail: thomson@cantab.net, jacob@bright.com

Course Times: 6:00pm-9:00pm, Tuesdays and Thursdays (Hattery)

Office Hours: Wednesday 5-7pm (preliminary)

# LOGISTICS

Course Website: https://www.schoology.com/

Instructors: Thomson Nguyen, Jacob Bollinger

E-mail: thomson@cantab.net, jacob@bright.com

Course Times: 6:00pm-9:00pm, Tuesdays and Thursdays (Hattery)

Office Hours: Wednesday 5-7pm (preliminary)

# I. WHAT IS DATA SCIENCE?
# II. THE DATA MINING WORKFLOW

# LAB:
# III. WORKING AT THE UNIX COMMAND LINE
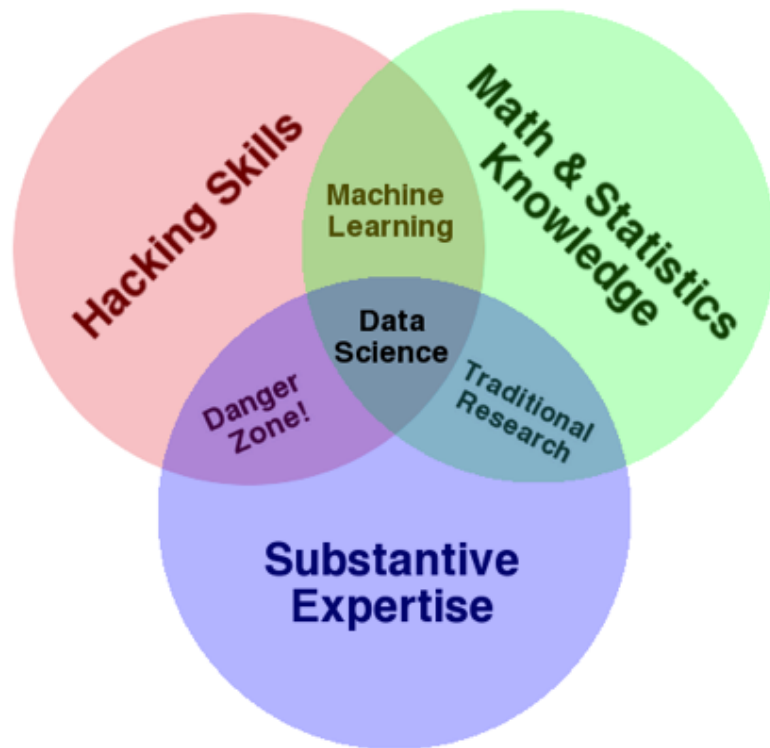# IV. VISUALIZING DATA WITH R & GGPLOT2

# I. WHAT IS DATA SCIENCE?

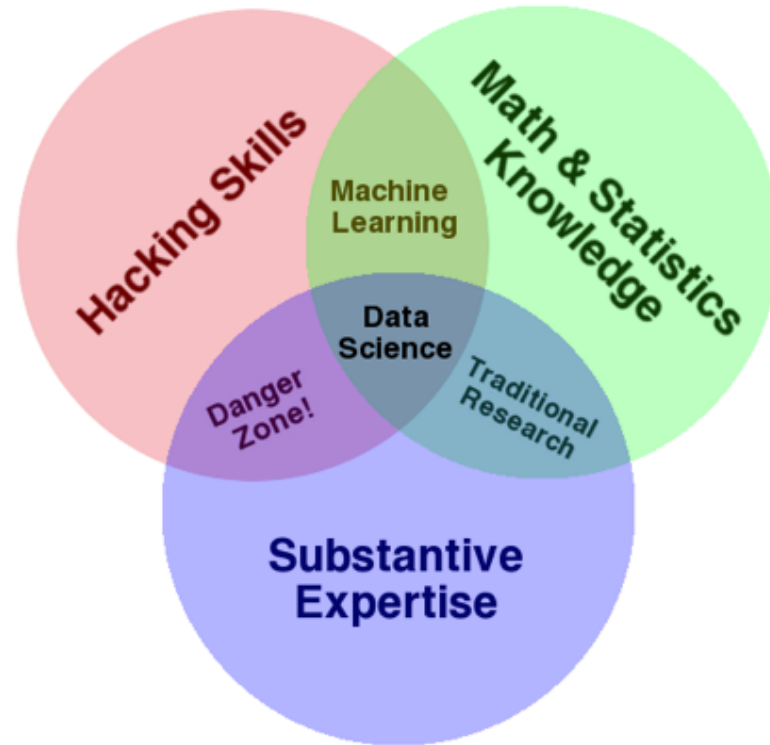‣ A set of tools and techniques used to extract useful information from data.

‣ A set of tools and techniques used to extract useful information from data.

‣ An interdisciplinary, problem-oriented subject.

source: http://www.dataists.com/2010/09/the-data-science-venn-diagram/

**ONE MORE THING!**

Communication skills

source: http://www.dataists.com/2010/09/the-data-science-venn-diagram/

‣ A set of tools and techniques used to extract useful information from data.

‣ An interdisciplinary, problem-solving oriented subject.

‣ The application of scientific techniques to practical problems.

‣ A set of tools and techniques used to extract useful information from data.

‣ An interdisciplinary, problem-solving oriented subject.

‣ The application of scientific techniques to practical problems.

‣ A rapidly growing field.

**Michael E. Driscoll**
@medriscoll

Following

Data scientists: better statisticians than most programmers & better programmers than most statisticians bit.ly/NHmRqu @peteskomoroch
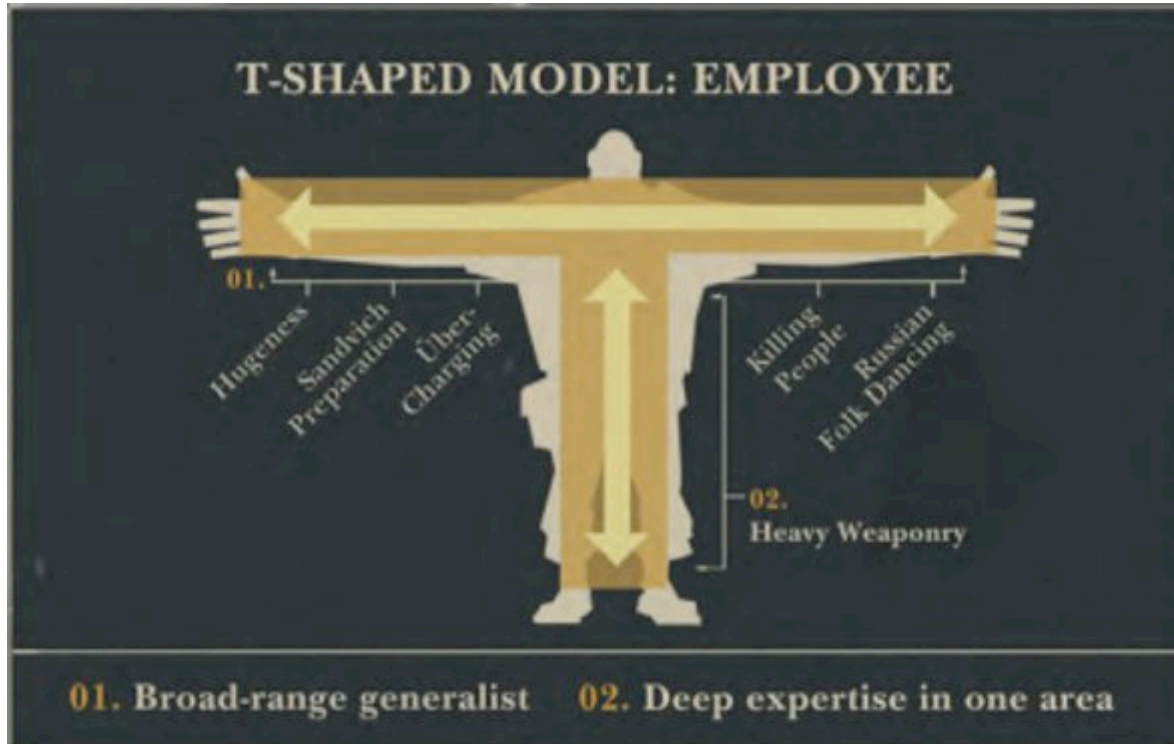
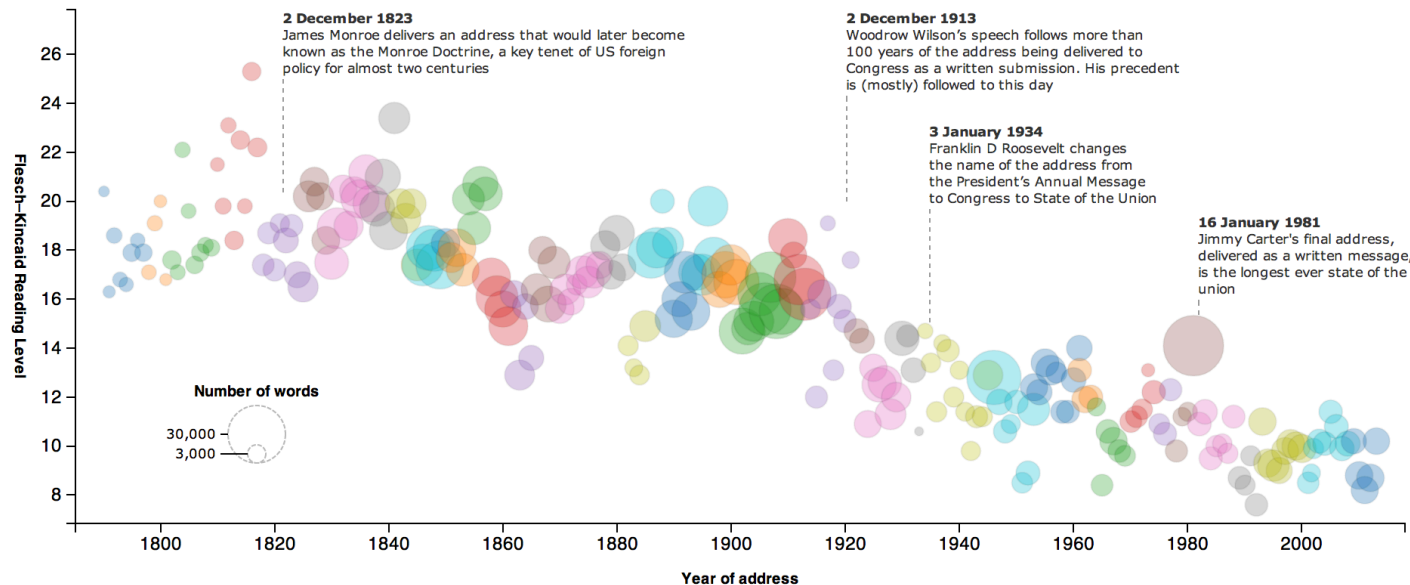← Reply  ⤺ Retweet  ★ Favorite  ••• More  ▼ Pocket

(Valve Software)

- Statistical and machine learning knowledge

- Engineering experience

- Academic curiosity

- Product sense

- Storytelling

- Cleverness

The state of our union is ... dumber:
How the linguistic standard of the presidential address has declined

Using the Flesch-Kincaid readability test the Guardian has tracked the reading level of every state of the union

# Music + Data:
# http://bit.ly/echonest

- Stack Overflow tag recommendation and response time prediction

- Locating ethnic food in ethnic neighborhoods

- Building optimal NBA teams

- Recommending new musical artists

- Prioritize emergency calls in Seattle

- Finding the right college for you

# II. THE DATA SCIENCE WORKFLOW

## Dataists

‣ 1. Obtain

‣ 2. Scrub

‣ 3. Explore

‣ 4. Model

‣ 5. Interpret

*source: http://berkeleydatascience.files.wordpress.com/2012/01/20120117berkeley1.pdf*

## Jeff Hammerbacher

- ‣ 1. Identify problem
- ‣ 2. Instrument data sources
- ‣ 3. Collect data
- ‣ 4. Prepare data (integrate, transform, clean, impute, filter, aggregate)
- ‣ 5. Build model
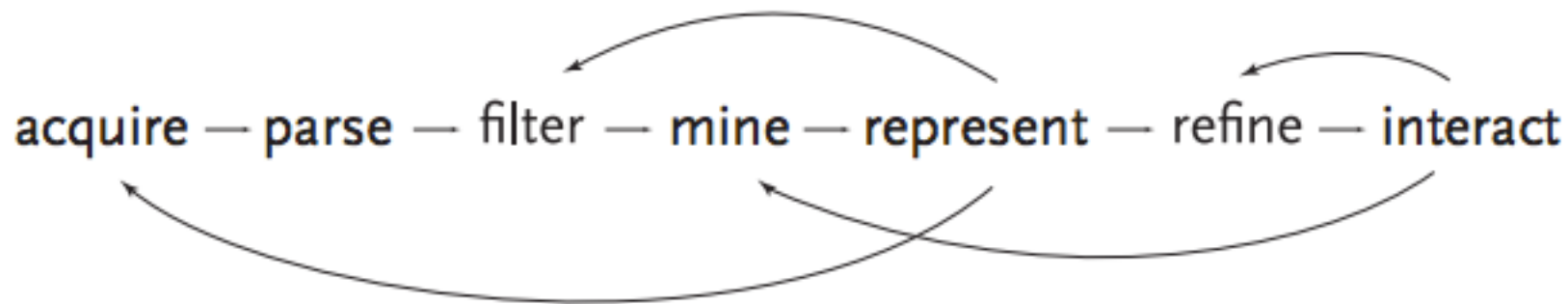- ‣ 6. Evaluate model
- ‣ 7. Communicate results

Ted Johnson

‣ 1. Assemble an accurate and relevant data set

‣ 2. Choose the appropriate algorithm

Ben Fry

‣ 1. Acquire
‣ 2. Parse
‣ 3. Filter
‣ 4. Mine
‣ 5. Represent
‣ 6. Refine
‣ 7. Interact

| COMPUTER SCIENCE | MATHEMATICS, STATISTICS, AND DATA MINING | GRAPHIC DESIGN | INFOVIS AND HCI |
|---|---|---|---|
| acquire    parse | filter    mine | represent    refine | interact |

source: http://benfry.com/phd/dissertation-110323c.pdf

acquire — parse — filter — mine — represent — refine — interact

**NOTE**

This diagram illustrates the *iterative* nature of problem solving

# III. WORKING AT THE UNIX COMMAND LINE

Download this dataset:

http://bit.ly/pacedataset

## KEY OBJECTIVES

- Navigate the filesystem
- Create, move, copy, and delete files & directories
- View & search files
- Edit & interact with files
- Combine steps
- Learn more

## TOOLS

- ls, cd
- cat, touch, mv, cp, mkdir, rm, rmdir

- head, tail, less, cat, grep
- vim, tr, sort, uniq, wc
- pipe (|)
- man, apropos

**NOTE**

Being comfortable at the command line makes your life much easier!

# IV. VISUALIZING DATA WITH R AND GGPLOT2

## KEY OBJECTIVES

- Become familiar with the R environment

- Explore data in R

- Visualize data using ggplot2

- Mathematical bonus: power laws

# IV. VISUALIZATIONS AS A MEDIUM

*Consider the following dataset:*
- *eleven (x, y) points*

*Consider the following dataset:*
- *eleven (x, y) points*
- *mean of x = 9, mean of y = 7.5*

Consider the following dataset:
- eleven (x, y) points
- mean of x = 9, mean of y = 7.5
- variance of x = 11, variance of y = 4.1

Consider the following dataset:
- eleven (x, y) points
- mean of x = 9, mean of y = 7.5
- variance of x = 11, variance of y = 4.1
- correlation of x and y = 0.8

*Consider the following dataset:*
- *eleven (x, y) points*
- *mean of x = 9, mean of y = 7.5*
- *variance of x = 11, variance of y = 4.1*
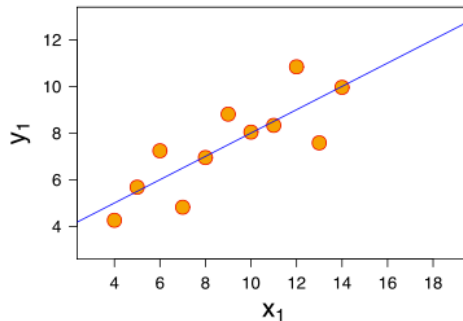- *correlation of x, y = 0.8*
- *line of best fit: y = 3.00 + 0.500x*

*Now, suppose I give you three more datasets with exactly the same characteristics...*

*Q: how similar are these datasets?*

*Now, suppose I give you three more datasets with exactly the same characteristics.*
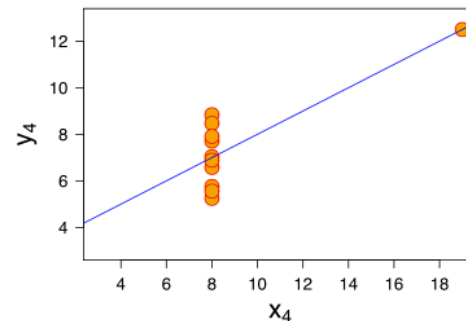
*Q: how similar are these datasets?*

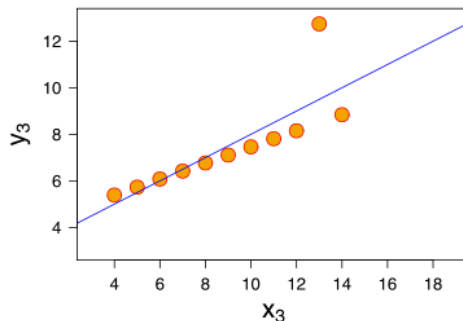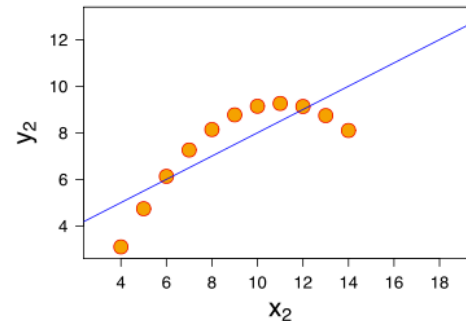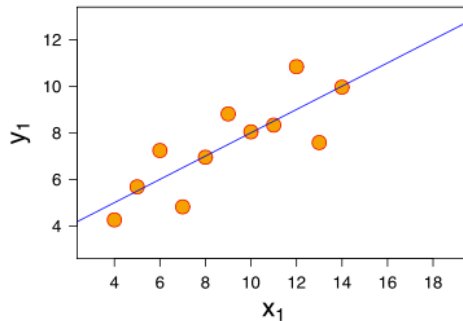*A: not very!*

http://en.wikipedia.org/wiki/Anscombe's_quartet

# DISCUSSION