# What is Data Science?

# Data Scientists: The Sexiest Job of the 21$^{st}$ Century
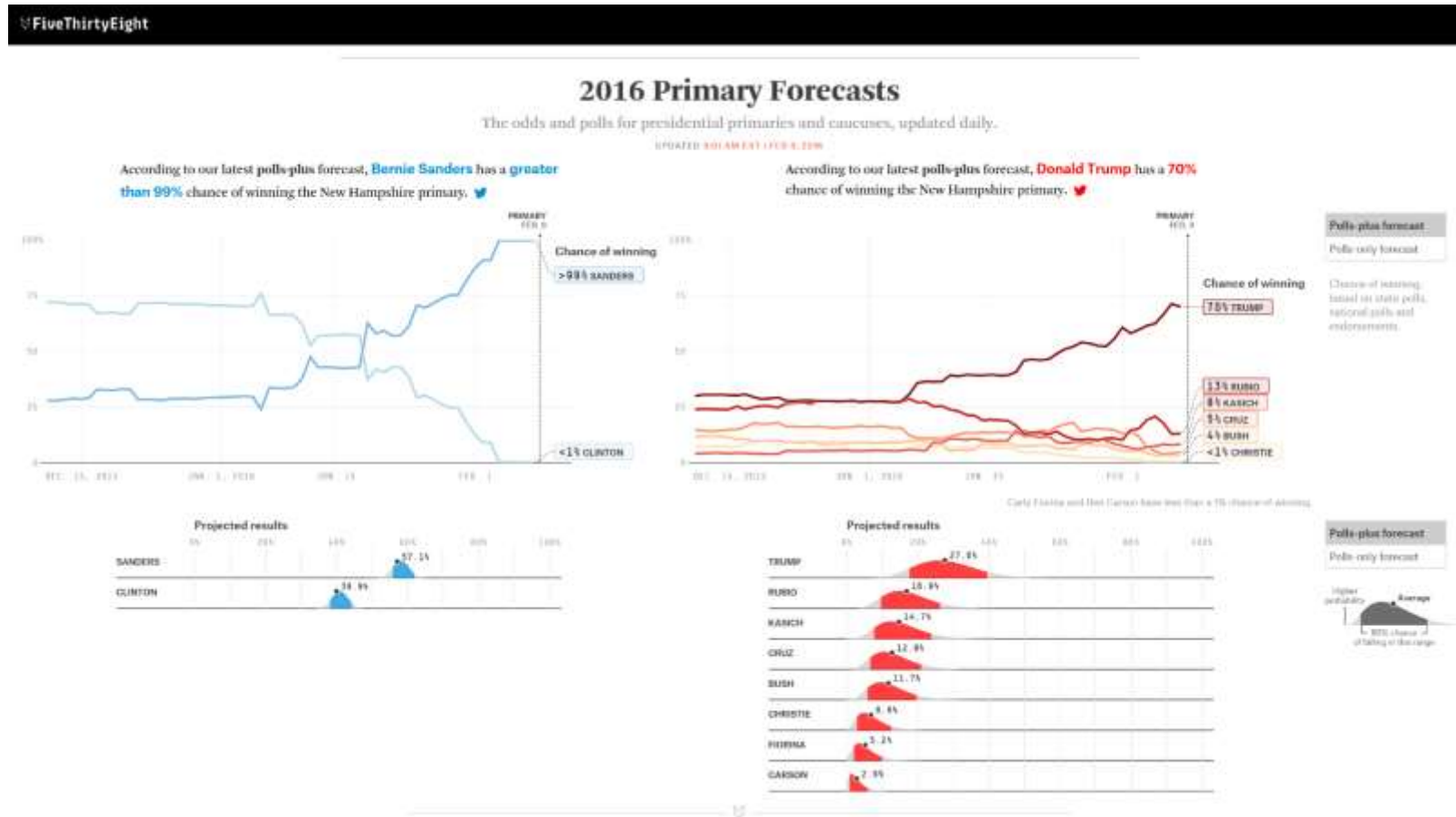




Source: Harvard Business Review

# FiveThirtyEight

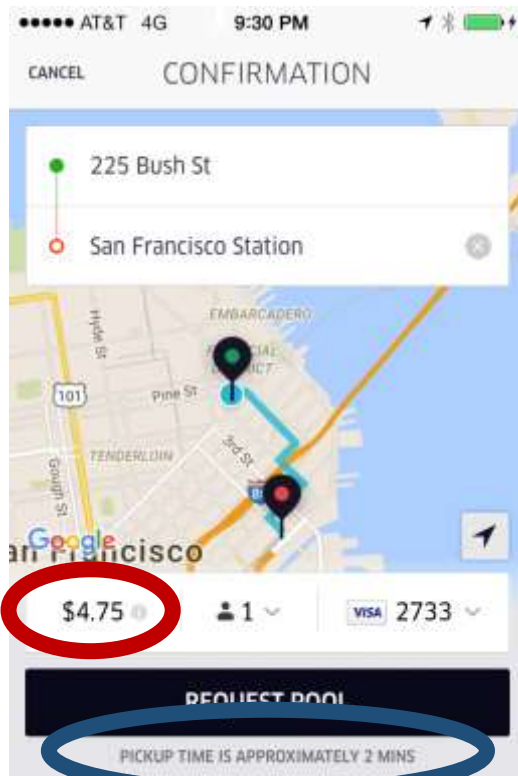# Data Science by FiveThirtyEight



Source: FiveThirtyEight

31

# Uber
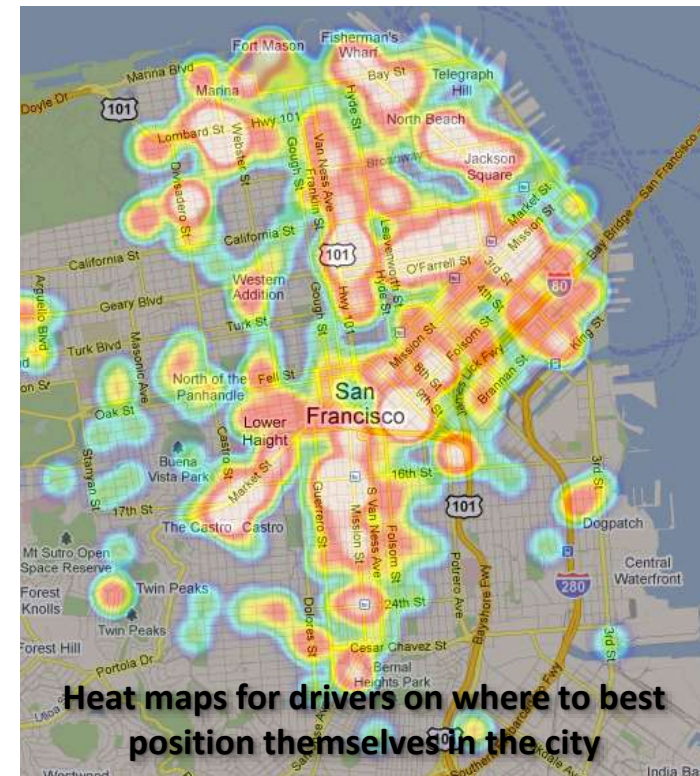
UBER

# Data Science by Uber

**For Riders...**

**and Drivers!**

UBER

**Fare Estimator and Dynamic Pricing (e.g., Surge Pricing) Algorithms**

**ETA (Estimated Time of Arrival) Algorithms**

**Heat maps for drivers on where to best position themselves in the city**

Sources: The Uber iPhone App/Wired

33

# Data Science Based-Business Models is the New Normal

Source: DOMO

35

# Who are Data Scientists?

# Activity: Who are Data Scientists?

**EXERCISE**

DIRECTIONS (10 minutes)

1. Who are Data Scientists?

2. How do Data Scientists add value?

3. What makes a good Data Scientist?

4. When finished, share your answers with your table

DELIVERABLE

Answers to the above questions

# Data Scientists in ≤140 characters

# What is Data Science? (cont.)

# Data is the New Oil of the Digital Economy and IoT, Big Data, DS, and Cloud relate to one another



Big Data

DS

IoT

Cloud

Kirill Kirsanov/Shao-Chun Wang © 123RF.com

# Data Science involves a variety of skillsets, not just one

|  | IDENTIFY the problem | ACQUIRE the data | PARSE the data | MINE the data | REFINE the data | BUILD a model | PRESENT the results |
|---|---|---|---|---|---|---|---|
| Computer Science / Software Engineering |  | X | X | X |  |  | X |
| Mathematics |  |  |  | X | X | X |  |
| Domain Expertise / Business | X |  | X |  |  |  | X |

Shao-Chun Wang © 123RF.com

# Data Science involves a variety of skillsets, not just one (cont.)



Source: Data Science for the C-suite

# Data Scientists have different roles that prioritize different skillsets but all roles involve some part of each skillset to form strong data teams

# To sum it up

‣ Data Science is an interdisciplinary field about processes and systems to extract knowledge or insights from data in various forms

‣ An (ideal) data scientist is "someone who has the both the engineering skills to acquire and manage large data sets, and also has the statistician's skills to extract value from the large data sets and present that data to a large audience" – John Rauser

# Data Science Workflow
## (and how it maps to the course)

# What is the Data Science Workflow for?

‣ A methodology for Data Science to produce *reliable* and *reproducible* results

   ‣ **Reliable**: Accurate findings

   ‣ **Reproducible**: Others can follow your steps and get the same results

‣ Similar to the scientific method

The scientific method:

   ‣ Ask a Question

   ‣ Do Background Research

   ‣ Construct a Hypothesis

   ‣ Test Your Hypothesis by Doing an Experiment

   ‣ Analyze Your Data and Draw a Conclusion

   ‣ Communicate Your Results

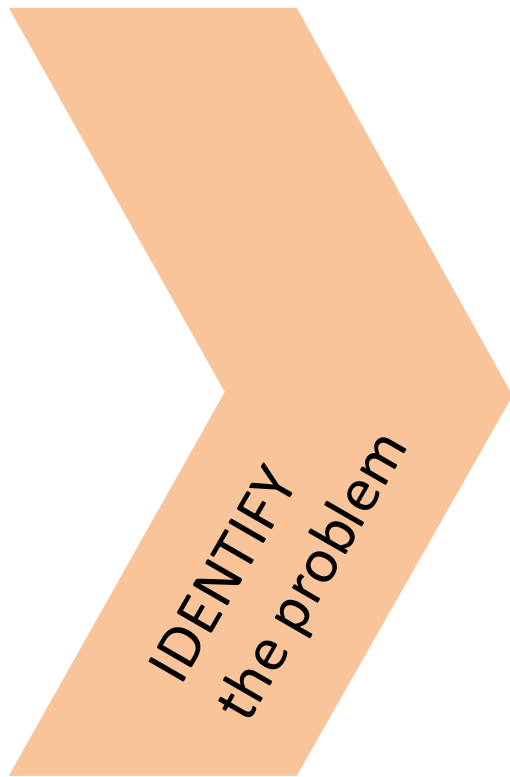# The Data Science Workflow
# (also called the Data Science Pipeline)



IDENTIFY the problem · ACQUIRE the data · PARSE the data · MINE the data · REFINE the data · BUILD a model · PRESENT the results
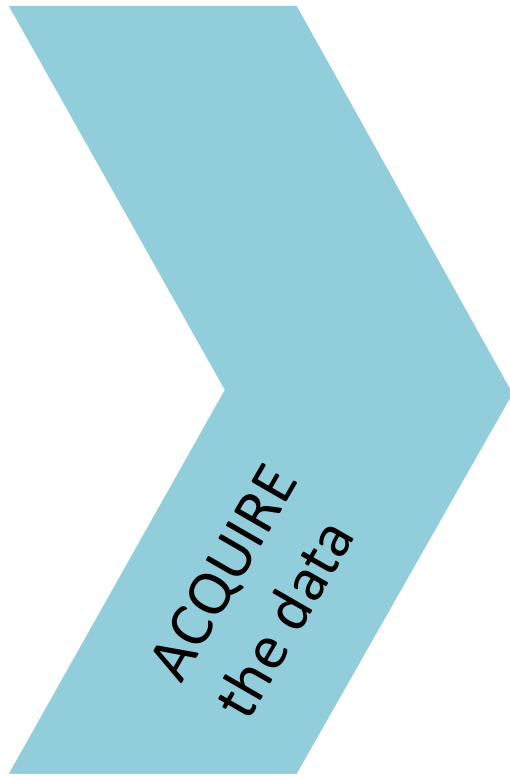
# ❶ Identify the Problem

IDENTIFY the problem

‣ Identify the Problem

‣ Identify business/product objectives

‣ Identify and hypothesize goals and criteria for success

‣ Create a set of questions for identifying correct dataset

# The Why's and How's of a Good Question



Corina Rosu © 123RF.com

# ❷ Acquire the Data
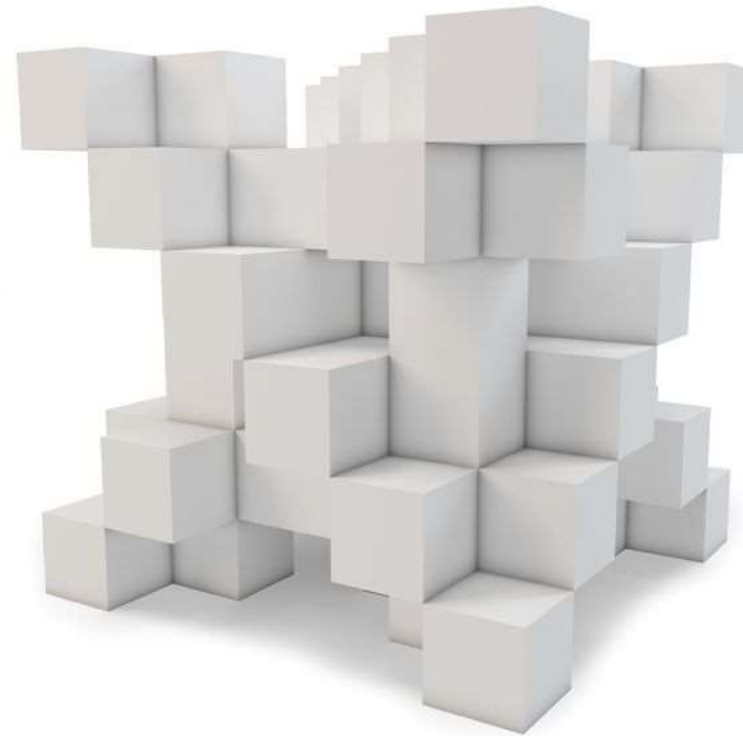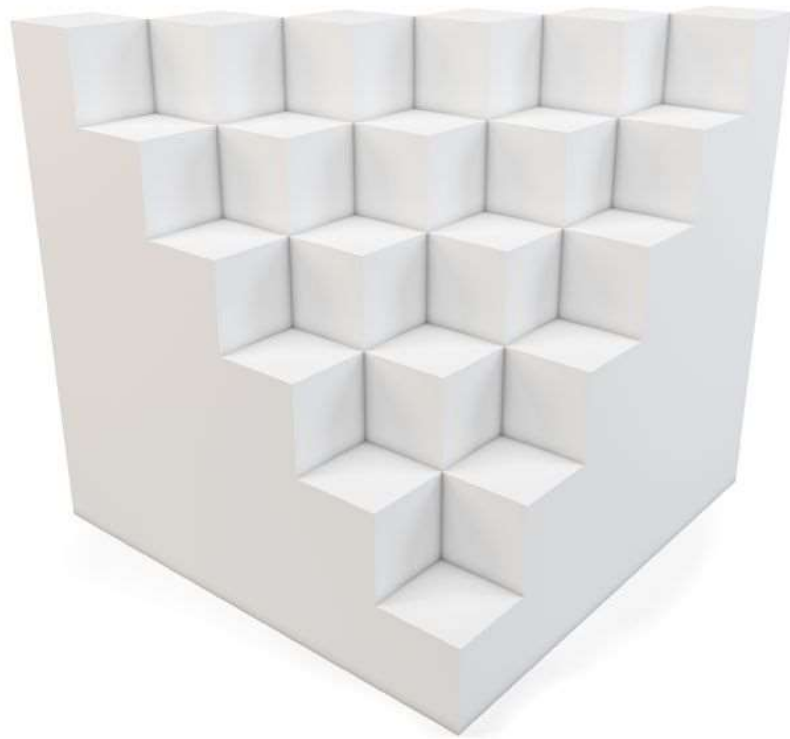
ACQUIRE the data

‣ Acquire the Data

  ‣ Identify the "right" dataset(s)

  ‣ Import data and set up local or remote data structure

  ‣ Determine most appropriate tools to work with data

# The data can be either unstructured or structured data

# What's an example of unstructured data?



Bundit Chuangboonsri © 123RF.com

‣ Sessions 13 and 14 in Unit 3

   ‣ Natural Language Processing

# However, most of the course will focus on structured data

‣ Unit 2

  ‣ Linear Regression (sessions 6 and 7)

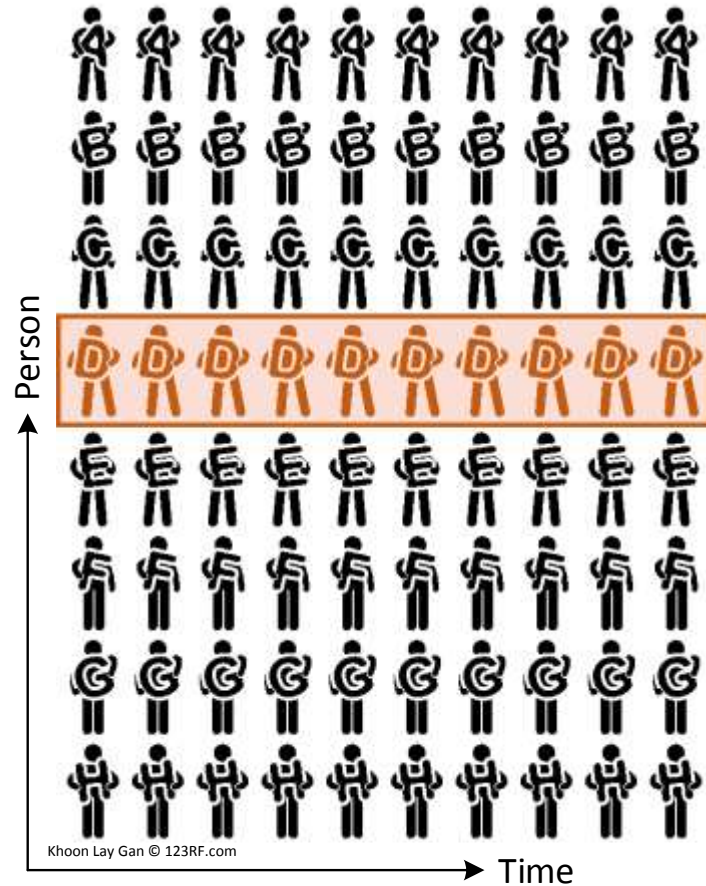  ‣ Classification and Logistic Regression (session 8 and 9)

‣ Unit 3

  ‣ Decision Trees and Random Forests (session 12)

milosb © 123RF.com

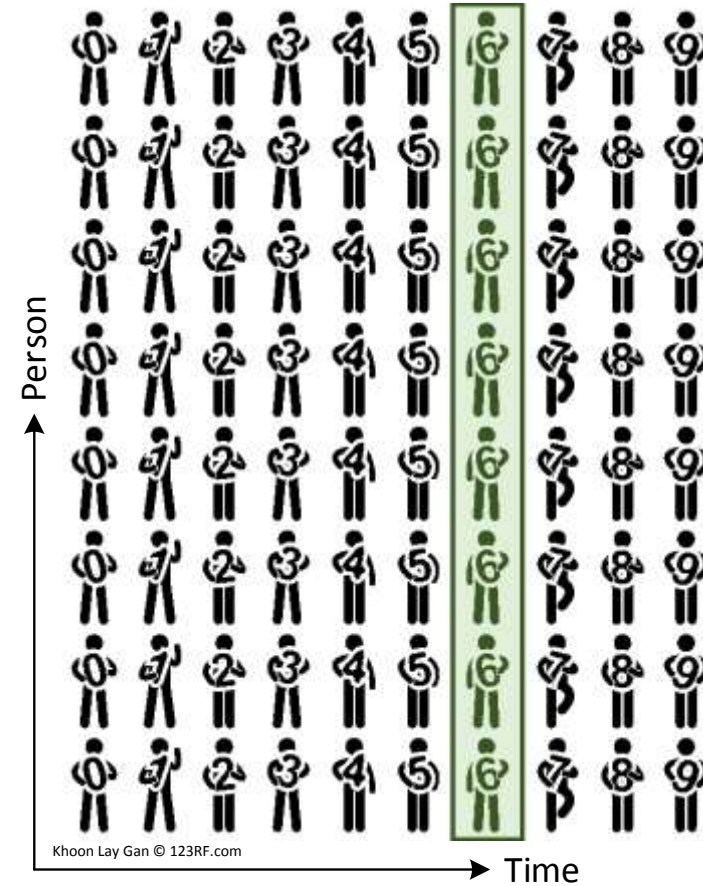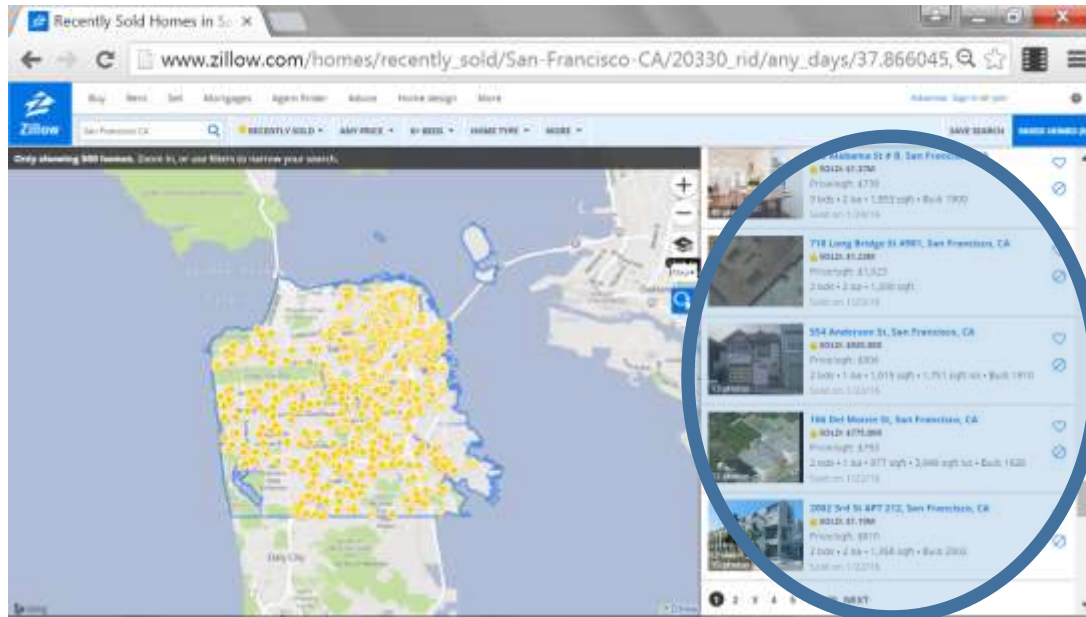# Unstructured data can be longitudinal



Person

Khoon Lay Gan © 123RF.com

Time

‣ Sessions 15 and 16 in Unit 3

    ‣ Time Series

# Unstructured data can be cross-sectional

‣ And most of the course will

focus on it



Person
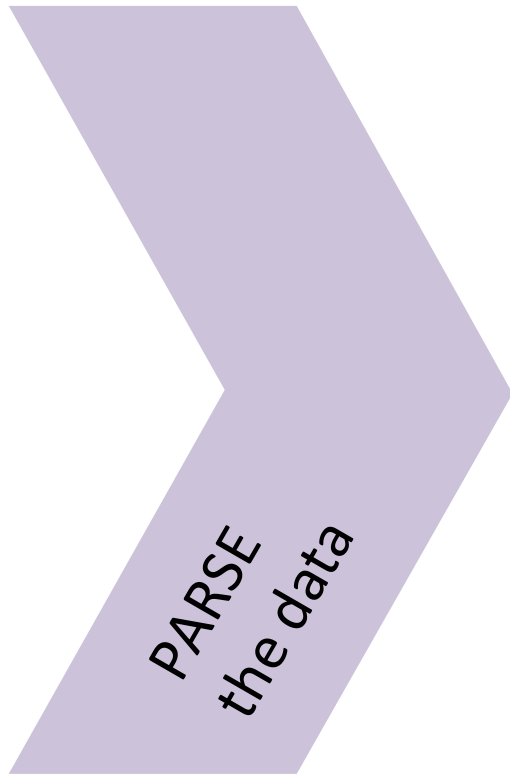
Time

Khoon Lay Gan © 123RF.com

# Raw structured data is Messy™…



```html
<div class="property-info"
id="yui_3_18_1_1_1456167242885_71870"><strong
id="yui_3_18_1_1_1456167242885_71869"><dt class="property-
address" id="yui_3_18_1_1_1456167242885_71868"><a
href="/homedetails/149-Shipley-St-San-Francisco-CA-
94107/15147894_zpid/" class="hdp-link routable" title="149
Shipley St, San Francisco, CA Real Estate"
id="yui_3_18_1_1_1456167242885_71873">149 Shipley St, San
Francisco, CA</a></dt></strong><dt class="listing-type zsg-
content_collapsed"
id="yui_3_18_1_1_1456167242885_71875"><span class="zsg-
icon-recently-sold type-icon"></span>Sold: $1.18M</dt><dt
class="zsg-fineprint"
id="yui_3_18_1_1_1456167242885_71877">Price/sqft:
$1,116</dt><dt class="property-data"
id="yui_3_18_1_1_1456167242885_71880"><span class="beds-
baths-sqft">3 bds • 2 ba • 1,057 sqft</span><span
class="built-year" id="yui_3_18_1_1_1456167242885_71879"> •
Built 1992</span></dt><dt class="sold-date zsg-fineprint"
id="yui_3_18_1_1_1456167242885_71975">Sold on
2/22/16</dt></div>
```
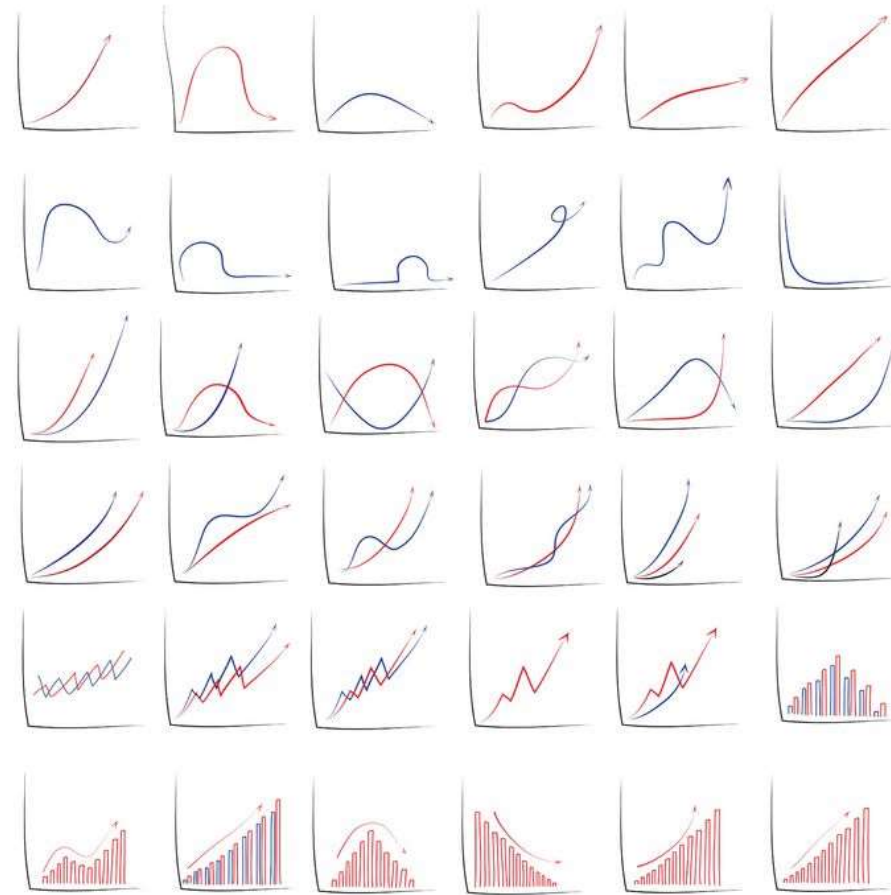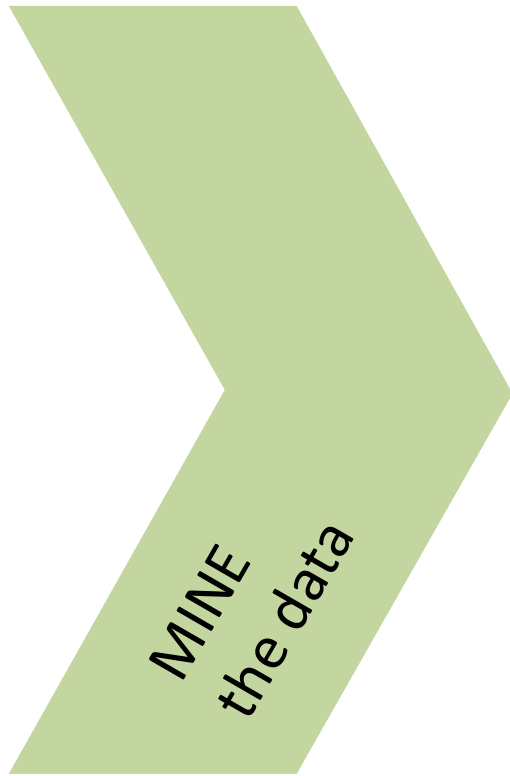
# ❸ Parse the Data

PARSE the data

‣ Parse the Data

   ‣ Read any documentation provided with the data

   ‣ Perform exploratory data analysis

   ‣ Verify the quality of the data

# Exploratory Data Analysis



Napat Polchoke © 123RF.com
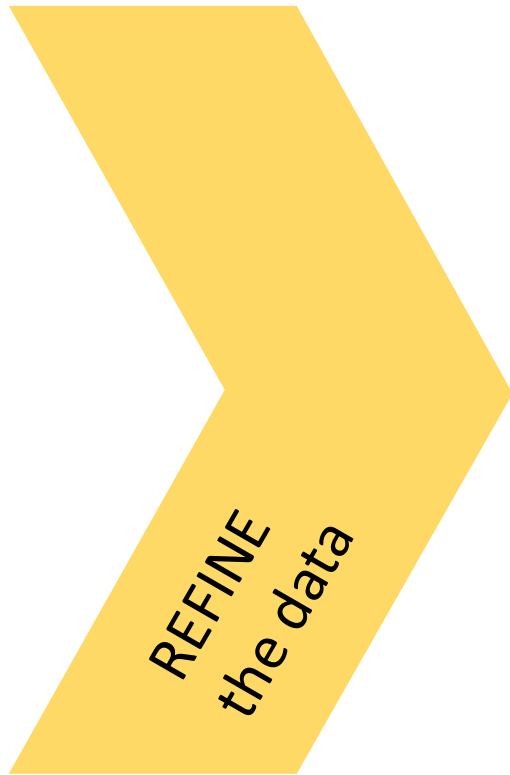
# ❹ Mine the Data

MINE the data

- ‣ Mine the Data

  - ‣ Determine sampling methodology and sample data

  - ‣ Format, clean, slice, and combine data in Python

  - ‣ Create necessary derived columns from the data (new data)

# We will be tidying our data using the Python *pandas* library

# ❺ Refine the Data

REFINE
the data

‣ Refine the Data

    ‣ Identify trends and outliers

    ‣ Apply descriptive and inferential statistics

    ‣ Document and transform data

# We will apply inferential statistics

# ❻ Build a Model

BUILD
a model

‣ Build a Model

‣ Select appropriate model

‣ Build model

‣ Evaluate and refine model

# The types of machine learning algorithms we will study in this course

| | Continuous | Categorical |
|---|---|---|
| **Supervised** (a.k.a., predictive modeling) | Linear Regression<br>K-Nearest Neighbors<br>Decision Trees and Random Forests | Logistic Regression<br>K-Nearest Neighbors<br>Decision Trees and Random Forests |
| *Unsupervised* | *A machine learning model that doesn't use labeled data is called unsupervised. It extract structure from the data. Goal is "representation"* | |

# ❼ Present the Results

PRESENT the results

‣ Present the Results

  ‣ Summarize findings with narrative, storytelling techniques

  ‣ Present limitations and assumptions of your analysis

  ‣ Identify follow up problems and questions for future analysis

# Know Your Audience



Corina Rosu © 123RF.com

# A Note About Iteration

‣ Iteration is an important part of *every* step in the Data Science Workflow.  At any given point in the process, you may find yourself repeating or going back and re-doing elements in order to better understand your data, clarify your model, and refine your presentation

‣ For example, after presenting your findings, you may want to:

   ‣ Identify follow-up problems and questions for future analysis

   ‣ Create a visually effective summary or report

   ‣ Consider the needs of different stakeholders and how your report might be changed for them

   ‣ Identify the limitations of your analysis

   ‣ Identify relationships between visualizations

# Multiple variants exist but they are pretty much all doing the same thing

- Jeff Hammerbacher
  - Identify problem
  - Instrument data sources
  - Collect data
  - Prepare data (integrate, transform, clean, impute, filter, aggregate)
  - Build model
  - Evaluate model
- Ben Fry
  - Acquire
  - Parse
  - Filter
  - Mine
  - Represent
  - Refine
  - Interact

- Peter Huber
  - Inspection
  - Error checking
  - Modification
  - Comparison
  - Modeling and model fitting
  - Simulation
  - What-if analyses
  - Interpretation
  - Presentation of conclusions
- Dataists
  - Obtain
  - Scrub
  - Explore
  - Model
  - Interpret

- Colin Mallows
  - Identify data to collect and its relevance to your problem
  - Statistical specification of the problem
  - Method selection
  - Analysis of method
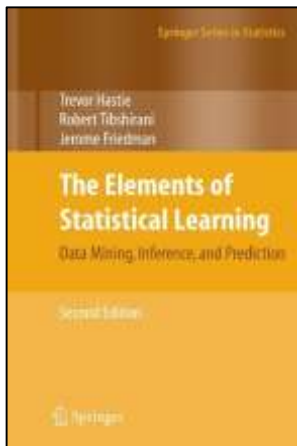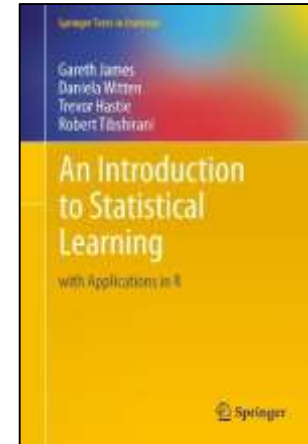  - Interpret results for non-statisticians
- Jim Gray
  - Capture
  - Curate
  - Communicate
- Ted Johnson
  - Assemble an accurate and relevant dataset
  - Choose the appropriate algorithm

# Some great resources to follow along the class (or afterwards) *(optional; not required for the course)*

‣ An Introduction to Statistical Learning: with Applications in R

(by James et al.).  The e-book is available free-of-charge [here](here)

‣ For a more advanced treatment of these topics, check out The Elements of Statistical Learning: Data Mining, Inference, and Prediction (by Hastie et al.).  And yes, the e-book is also free… ([here](here))
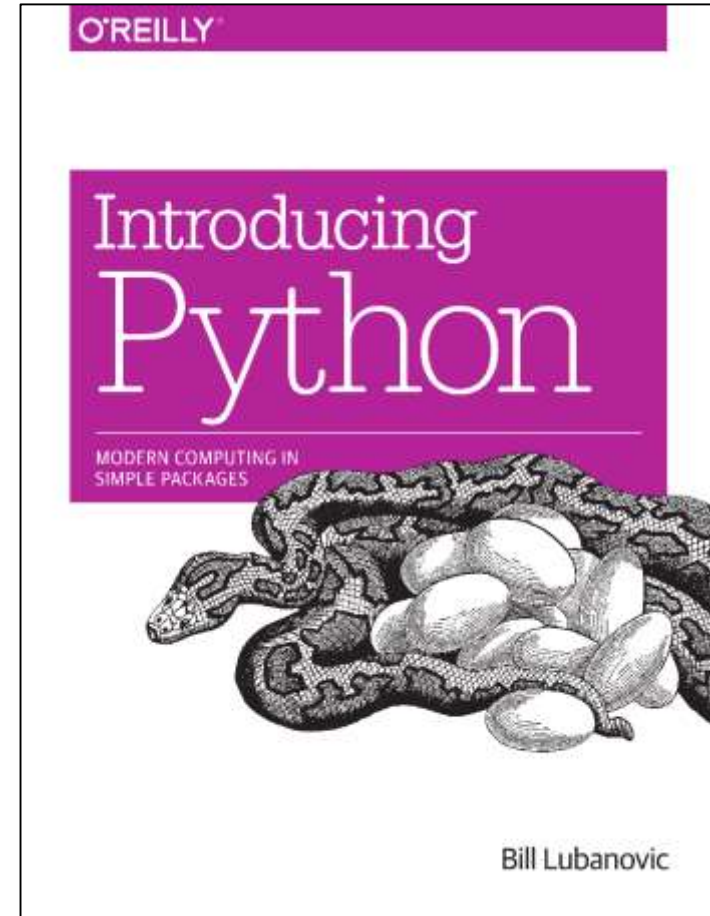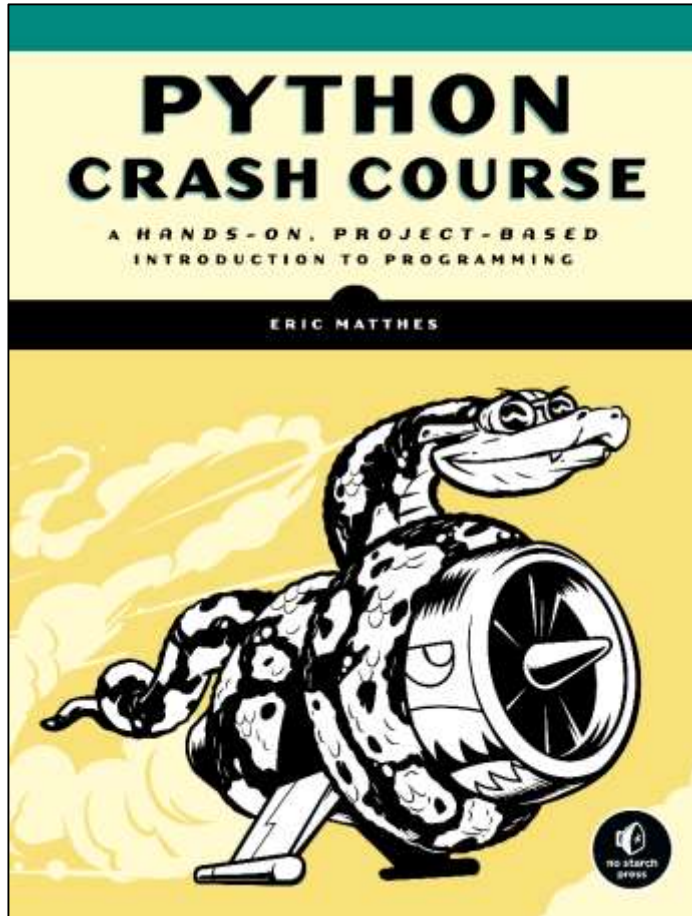
# Q & A

# Onboarding Review

# Python and iPython Notebooks

A couple of resources to get started with Python *(optional; not required for the course)* (the contents overlap but the styles are different so you should only pick one if any)

# Review

# Review

You should now be able to answer the following questions:

‣ What is data science?

‣ What is the data science workflow?

‣ How can you have a successful learning experience at GA?

# Q & A

# Next Class

*Research Design and pandas*

# Learning Objectives

After the next lesson, you should be able to:

‣ Setup and manage your personal GitHub repository for submitting assignments

‣ Define a problem and types of data

‣ Identify dataset types

‣ Apply the data science workflow in the *pandas* context

‣ Write an iPython notebook to import, format, and clean data using the *pandas* library

# Exit Ticket

*Don't forget to fill out your exit ticket here*