

What is Data Science

Ivan Corneillet

Data Scientist

Learning Objectives

After this lesson, you should be able to:

- Describe the components of a successful learning environment
- Define what is data science and who data scientists are
- Setup your development environment and learn the different workflows we will use in this course
- Define the data science workflow
- Define a data science problem
- Identify the components of a concise and convincing report and how they relate to specific audiences/stakeholders

Here's what's happening today:

- Welcome to GA and DS!
- Setting you up for success
- What is Data Science
 - and who are data scientists?
- Class Discussion
 - “Data Scientists: The Sexiest Job of the 21st Century” (Harvard Business Review)
 - “The Rise of Artificial Intelligence and the End of Code” (Wired)
- Installfest
- The Data Science Workflow
 - Overview
 - ❶ IDENTIFY the Problem
 - ❷ PRESENT the Results
- Exit Tickets
- Homework – Python Review

A black circle containing the white text "DS".

DS

Welcome to GA and DS!

A black circle containing the white letters "DS" in a bold, sans-serif font.

DS

Setting You Up for Success

Meet Your Team

- Ivan Corneillet, Lead Instructor



- Megan O'Rourke, Associate Instructor

- Matt Jones, Course Producer



Course Logistics

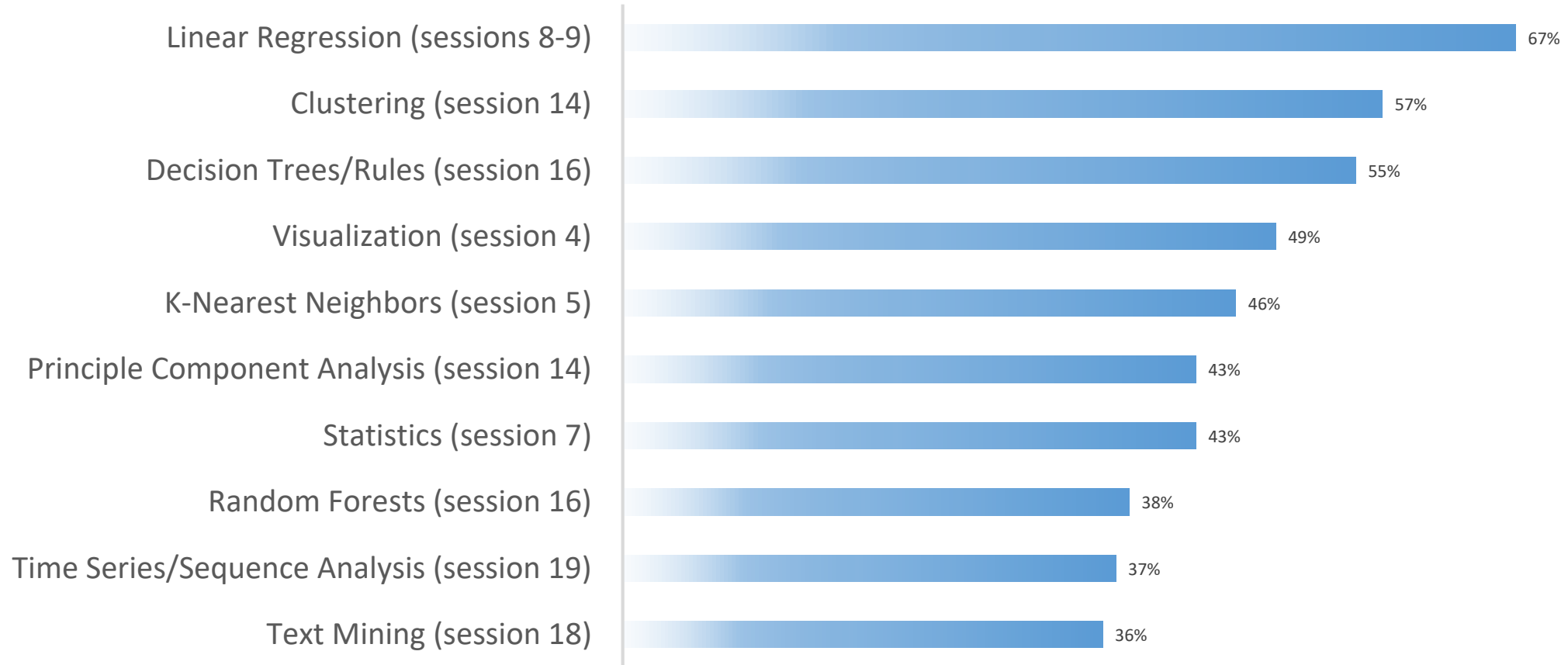
- Lead Instructor
 - Ivan Corneillet (ivan@paspieur.com)
- Associate Instructor
 - Megan O'Rorke (megororke@gmail.com)
- Course Producer
 - Matt Jones (studentservicesSF@ga.co)
- Class
 - December 6 – February 21, Tuesdays and Thursdays, 6:30PM – 9:30PM (no classes 12/19–12/31)
 - Classroom 5
- Slack
 - <https://ds-sf-30.slack.com>
- GitHub
 - <https://github.com/ga-students/DS-SF-30>
- Exit Tickets
 - <http://tiny.cc/ds-sf-30>

What skills will I learn in this class?

What is Data Science <i>(session 1)</i>	Research Design <i>(session 1)</i>	The <i>pandas</i> library <i>(sessions 2, 4, and 6)</i>	Databases and APIs <i>(session 3)</i>	Exploratory Data Analysis (and Data Visualization) <i>(sessions 4 and 6)</i>
k-Nearest Neighbors <i>(sessions 5, 12, and 17)</i>	Model Fit <i>(sessions 5, 7–9, 12–13, and 17)</i>	Linear Regression <i>(sessions 8–9, 12, and 17)</i>	Regularization <i>(sessions 10, 12, and 17)</i>	Logistic Regression <i>(sessions 11–12, and 17)</i>
Clustering <i>(sessions 14 and 17)</i>	Trees <i>(sessions 16–17)</i>	Natural Language Processing <i>(session 18)</i>	Time Series <i>(session 19)</i>	Presenting Insights from Data Models <i>(sessions 1, 15, and 20)</i>

Top 10 Algorithms & Methods used by Data Scientists

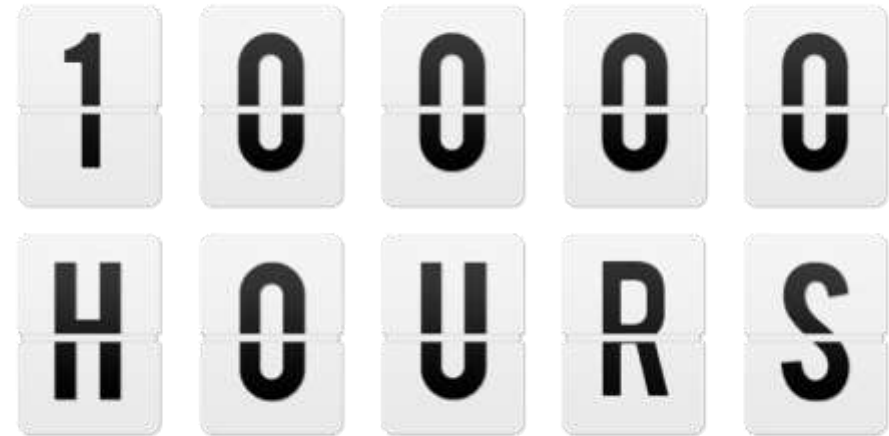
(<http://www.kdnuggets.com/2016/09/poll-algorithms-used-data-scientists.html>)



Gladwell's 10,000 Hour Rule

(<http://www.wisdomgroup.com/blog/10000-hours-of-practice>)

- ▶ “Greatness requires enormous time”
 - ▶ It takes roughly ten thousand hours of practice to achieve mastery in a field



How will I apply and reinforce these new skills?

Unit Project You will design a research project, perform exploratory data analysis and build a logistic model to determine what factors affect admission the most	Research Design		Exploratory Data Analysis		Machine Learning Modeling		Executive Summary			
Final Project Using a dataset of your choosing, you will design a project, build a data science model and present their finding to the course	Lightning Presentation		Experimental Write-up		Exploratory Data Analysis		Notebook Draft		Project Presentation	

Typical Class

- *Pre-readings (usually optional)*
- Today's objectives
- Announcements and exit tickets
- Review of the previous class
- Series alternating between:
 - Lectures
 - (deck, whiteboard, codealongs, and demos)
 - Practices
 - (cold-calling, individual and group exercises, and codealongs)
- Review of today's class
- Exit tickets
- Homework (*ungraded*)
- *Post-readings (usually optional)*



DS

What is Data Science and Who are Data Scientists?

Harvard Business Review | “Data Scientists: The Sexiest Job of the 21st Century” (2012)

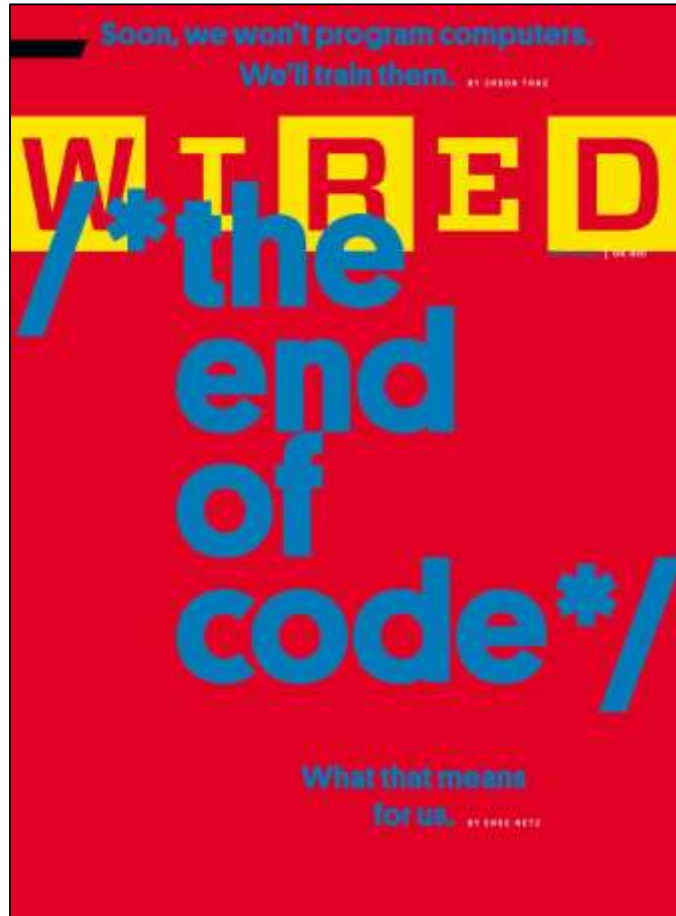
(<https://hbr.org/2012/10/data-scientist-the-sexiest-job-of-the-21st-century/>)



Source: Harvard Business Review

Wired | “The Rise of Artificial Intelligence and the End of Code” (2016)

(<http://www.wired.com/2016/05/the-end-of-code/> and <https://www.wired.com/2016/05/google-alpha-go-ai/>)



Source: Wired

Activity | What is Data Science and Who are Data Scientists?



EXERCISE


DIRECTIONS (15 minutes)

1. What is data science? What are its applications? Why now? What's next?
2. Who are data scientists? How do they add value? What makes a good data scientist?
3. When finished, share your answers with your table

DELIVERABLE

Answers to the above questions

Data science is everywhere

 **FiveThirtyEight**

NETFLIX



Walmart 

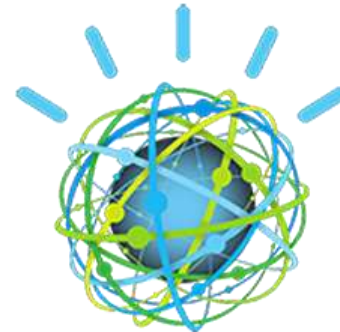


amazon 

Google



U B E R



IBM **Watson**

Linked in

Common questions asked in data science

How much? How many?

- What will the temperature be next Tuesday?
- What will my fourth quarter sales in France be?
- How many kilowatts will be demanded from my wind farm 30 minutes from now?
- How many new followers will I get next week?

Regression

- Predict a continuous outcome
 - k-Nearest Neighbors (session 5)
 - Linear Regression (sessions 8 and 9)
 - Trees (session 16)

Common questions asked in data science (cont.)

Is this A, B or C?

- Will this customer default on their loan?
- Is this an image of a man, a cat, or a dog?
- Will this customer click on the advertisement?
- Which team will win the championship?
- Is this mole malignant or benign?

Classification

- Predict a discrete outcome
 - k-Nearest Neighbors (session 5)
 - Logistic Regression (session 11)
 - Trees (session 16)

Common questions asked in data science (cont.)

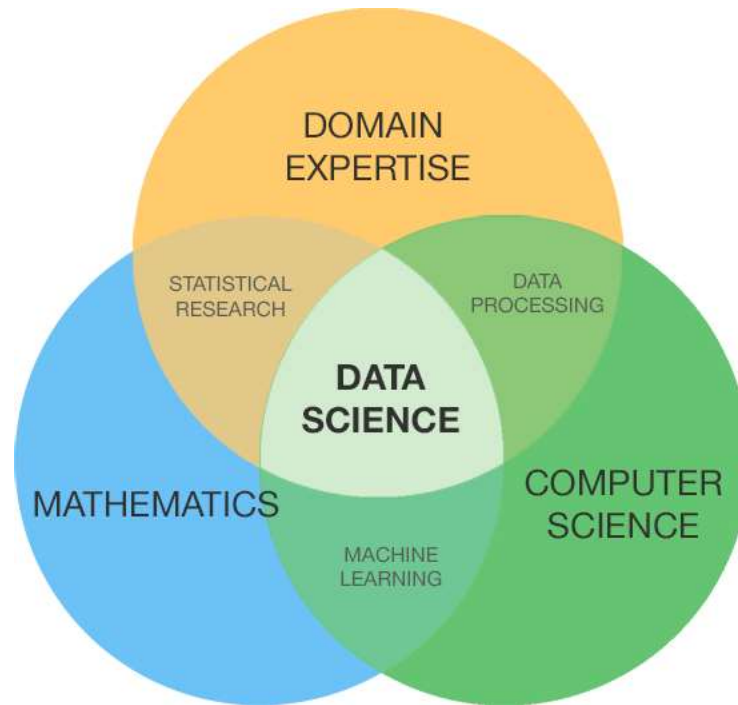
How is this Data Organized?

- What are the different types of coffee drinkers?
- Which viewers like the same kind of movies?
- What kinds of car models does GM produce?
- Are there common clusters of cable channels that customers tend to purchase together
- What is a natural way to break these documents into five topics?

Clustering

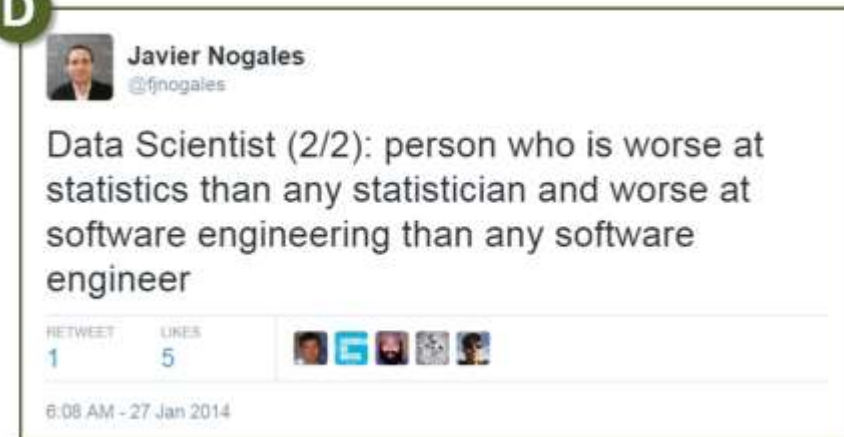
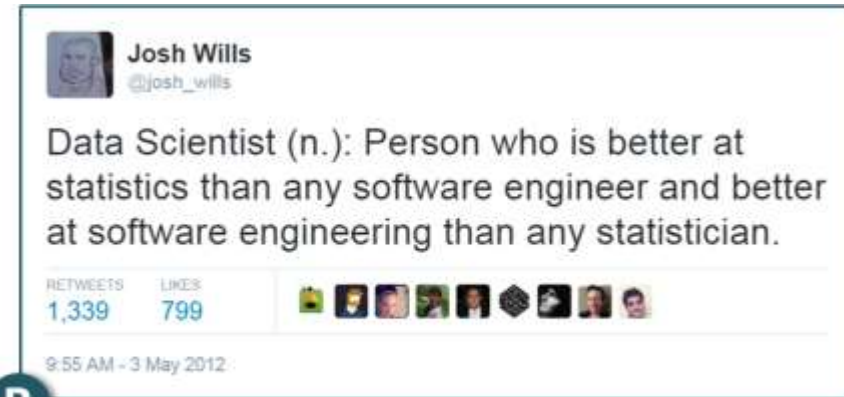
- What are the “categories” within the data?

Data science involves a variety of skillsets



Source: Data Science for the C-suite

Data scientists in ≤ 140 characters



Source: Twitter

Wired's “The Rise of Artificial Intelligence and the End of Code” (2016)

Behaviorism/Behavioral Psychology

- Brain as a black box
 - Stimulus and response, feedback and reinforcements
 - “ring bell, dog salivates”

Cognitive Psychology

- Brain more like a computer
 - Thoughts as programs
 - Absorb, process, and act upon information

Wired's "The Rise of Artificial Intelligence and the End of Code" (2016)

Machine Learning

- Humans *train* computers
 - Keep showing cats to a computer and eventually it will *learn* to recognize cats (<https://www.wired.com/2012/06/google-x-neural-network/>)
 - No symbols, no rules; instead an unparsable machine learning

Traditional Programming

- Humans *write code* (as explicit step-by-step-instructions) for computers to follow
 - Rule-based determinism
 - "Write enough rules and eventually, we'd create a system sophisticated enough to understand the world"
 - For years, Google Search relied mostly on these human-written rules (<https://www.wired.com/2016/02/ai-is-changing-the-technology-behind-google-searches/>)

Wired's "The Rise of Artificial Intelligence and the End of Code" (2016)

Age of Entanglement

- Outside-in view of how machine work
 - "Code doesn't just determine behavior, behavior also determine code"

Age of Enlightenment

- Inside-out view of how machine work
 - "First, we write the code, then the machine expresses it"

In this class we will model the stimuli as a matrix X (the feature matrix); the response is modelled as a vector y (the response vector). These are very important data structures that will be used as inputs to our machine learning algorithms

Feature Matrix X

Stimulus/feedback
"ring bell"

	col0	col1	col2	col3
row0				
row1				
row2				
row3				

Response Vector y

Response/reinforcements
"dog salivates"

	col
row0	
row1	
row2	
row3	



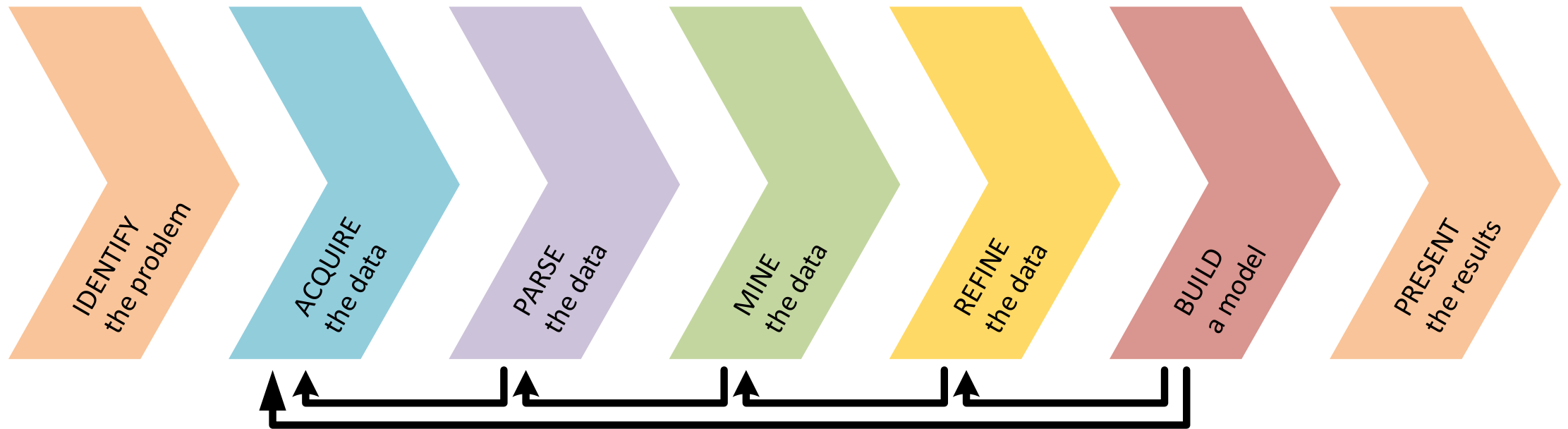
Installfest



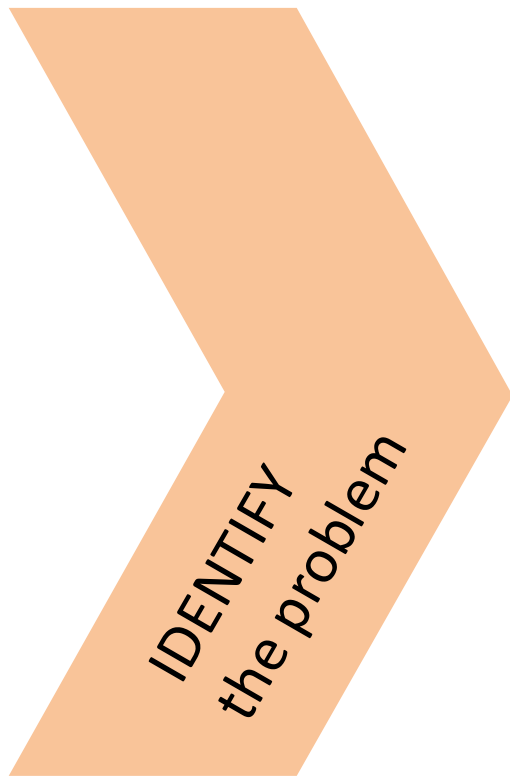
DS

Data Science Workflow

The Data Science Workflow



① 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

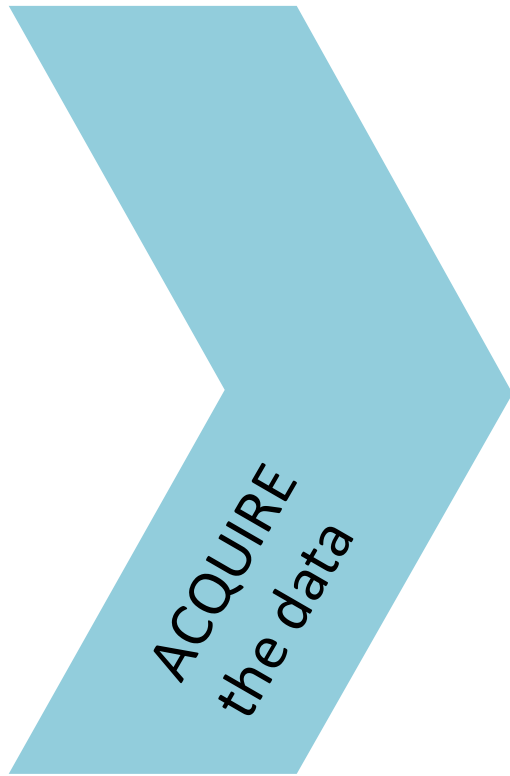
① IDENTIFY the Problem

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



Corina Rosu © 123RF.com

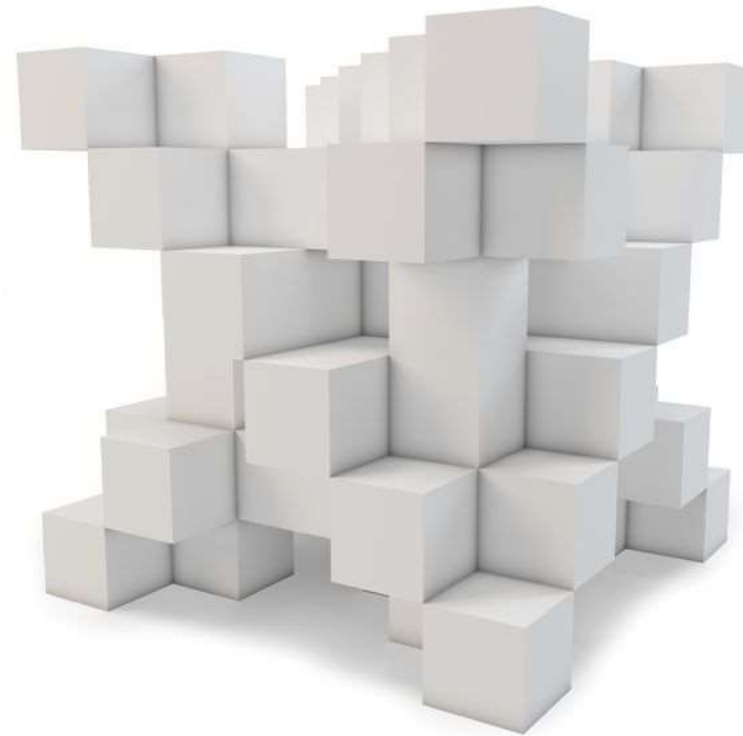
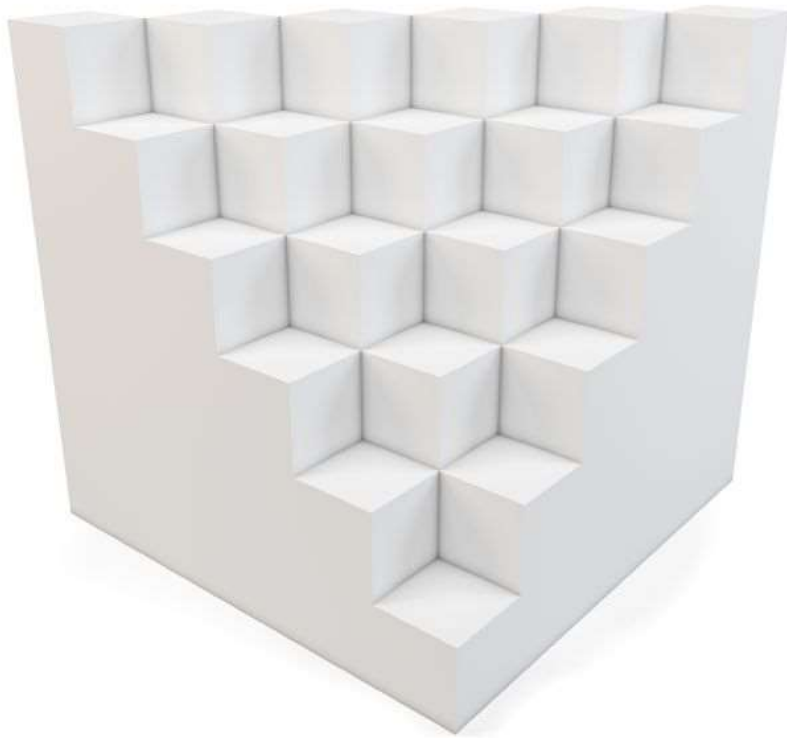
② ACQUIRE the Data



- 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

② ACQUIRE the Data

Data can be either unstructured or structured data



② ACQUIRE the Data

What's an example of unstructured data?

- Natural Language Processing
(session 18)



Bundit Chuangboonsri © 123RF.com

② ACQUIRE the Data

Most of the course will focus on structured data

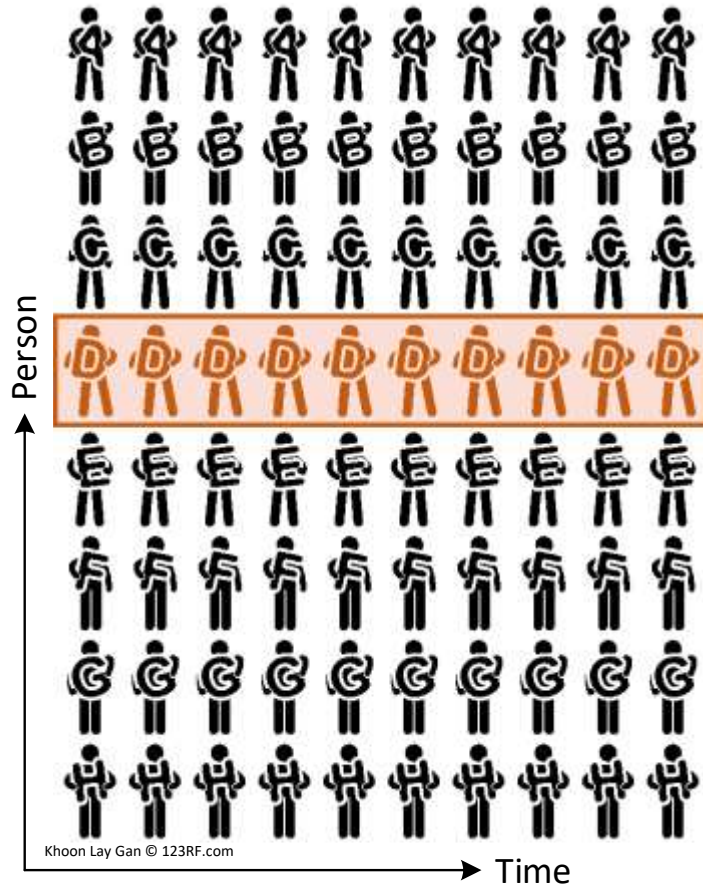
- k-Nearest Neighbors (*session 5*)
- Linear Regression (*sessions 8 and 9*)
- Logistic Regression (*session 11*)
- Trees (*session 16*)



milosb © 123RF.com

② ACQUIRE the Data

Unstructured data can be longitudinal

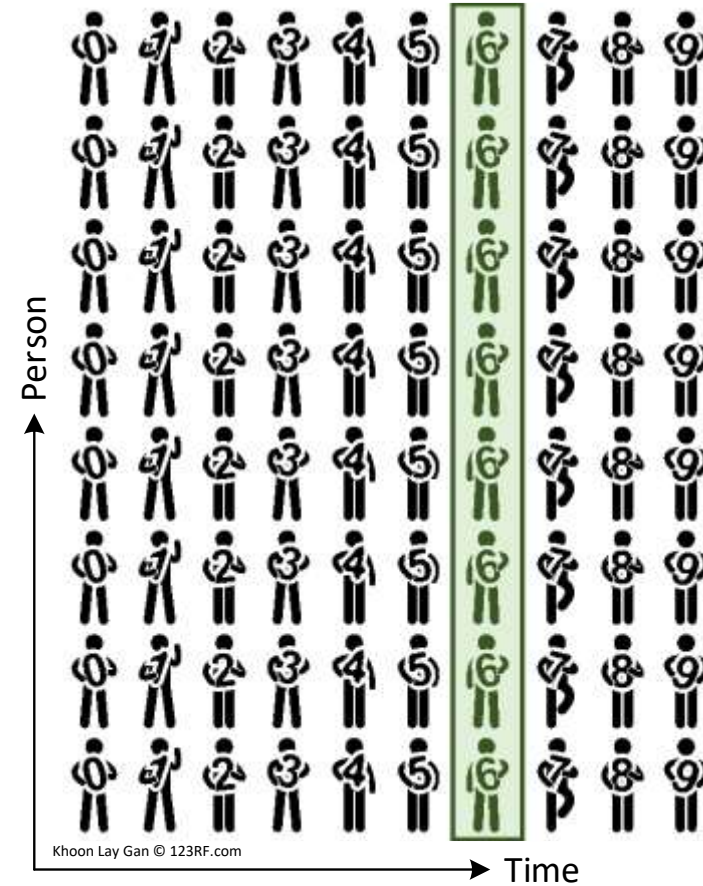


▸ Time Series (*session 19*)

② ACQUIRE the Data

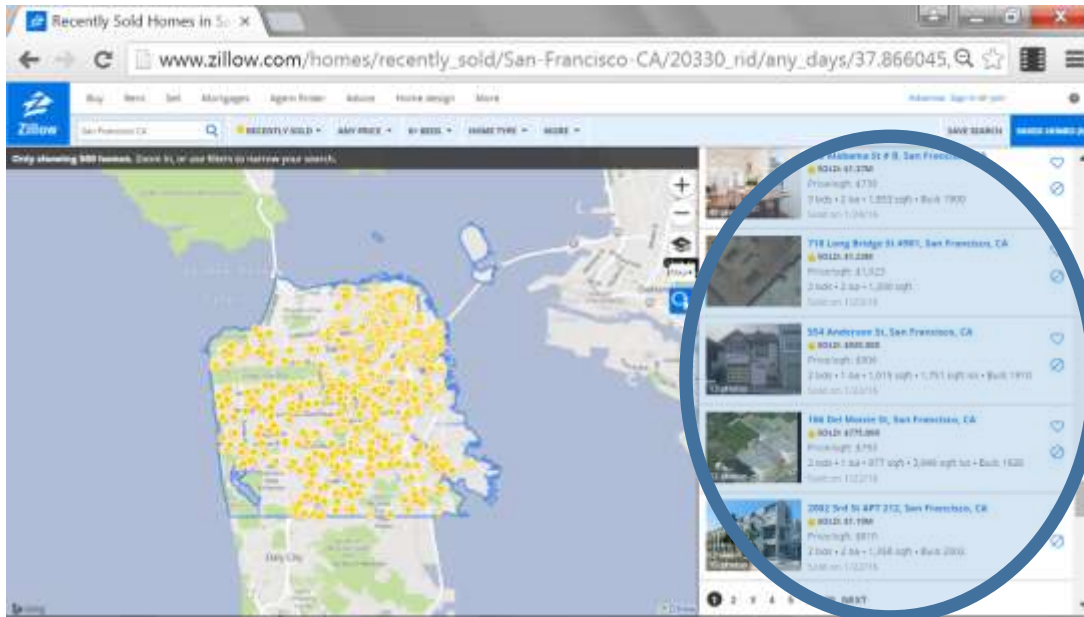
Unstructured data can be cross-sectional

- And most of the course will focus on it



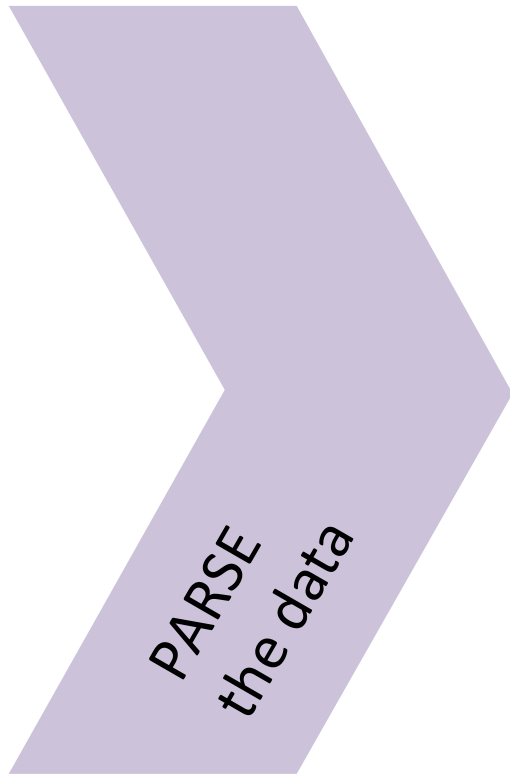
② ACQUIRE the Data

Raw structured data is Messy™...



```
<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>
```

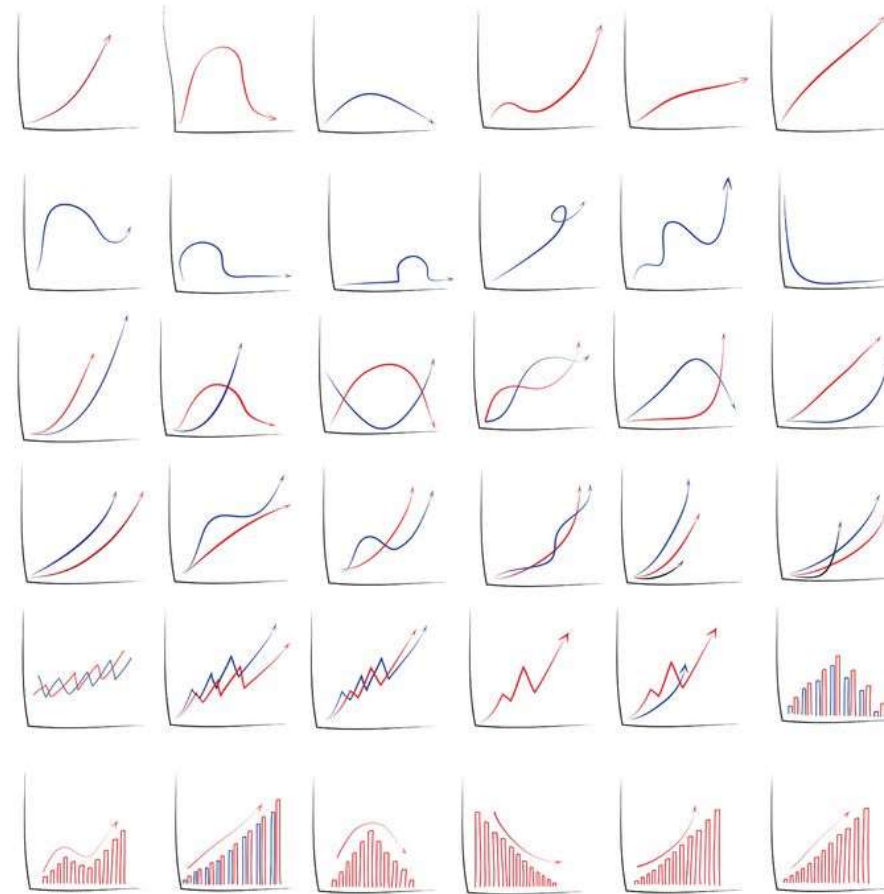
③ PARSE the Data



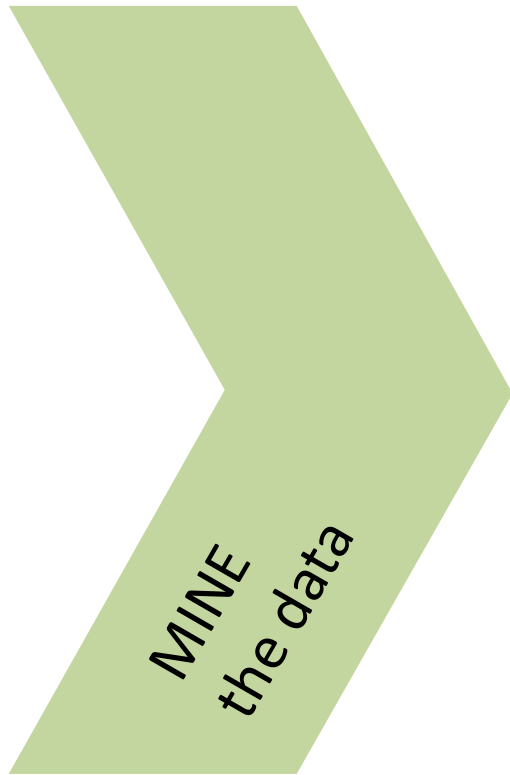
- Parse the Data
 - Read any documentation provided with the data
 - Perform exploratory data analysis
 - Verify the quality of the data

③ PARSE the Data

Exploratory Data Analysis



④ 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)

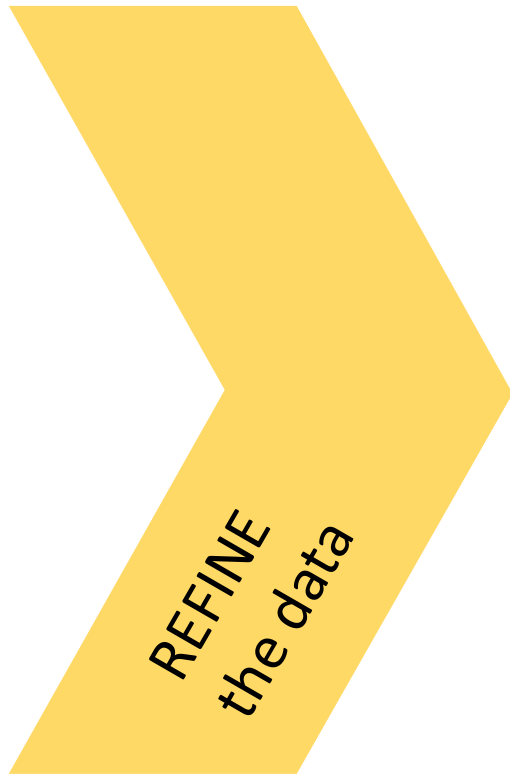
④ MINE the Data

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

The screenshot displays an Excel spreadsheet with the following data:

	A	B	C	D	E	F	G	H	I	J	K	L	M
	ID	Address	Latitude	Longitude	DateOfSale	SalePrice	SalePriceUnit	IsAStudio	BedCount	BathCount	Size	SizeUnit	Location
2	15063340	44444444	37799474	-122414835	11/30/2015	1.29 \$M		FALSE	2	2	1165 sqft		N/A
3	15063340	44444444	37799474	-122414835	11/30/2015	1.29 \$M		FALSE	2	2	1165 sqft		N/A
4	15063340	44444444	37799474	-122414835	11/30/2015	1.29 \$M		FALSE	2	2	1165 sqft		N/A
5	15063340	44444444	37799474	-122414835	11/30/2015	1.29 \$M		FALSE	2	2	1165 sqft		N/A
6	15063340	44444444	37799474	-122414835	11/30/2015	1.29 \$M		FALSE	2	2	1165 sqft		N/A
7	15063340	44444444	37799474	-122414835	11/30/2015	1.29 \$M		FALSE	2	2	1165 sqft		N/A
8	15063340	44444444	37799474	-122414835	11/30/2015	1.29 \$M		FALSE	2	2	1165 sqft		N/A
9	15063340	44444444	37799474	-122414835	11/30/2015	1.29 \$M		FALSE	2	2	1165 sqft		N/A
10	15063340	44444444	37799474	-122414835	11/30/2015	1.29 \$M		FALSE	2	2	1165 sqft		N/A
11	15063340	44444444	37799474	-122414835	11/30/2015	1.29 \$M		FALSE	2	2	1165 sqft		N/A
12	15063340	44444444	37799474	-122414835	11/30/2015	1.29 \$M		FALSE	2	2	1165 sqft		N/A
13	15063340	44444444	37799474	-122414835	11/30/2015	1.29 \$M		FALSE	2	2	1165 sqft		N/A
14	15063340	44444444	37799474	-122414835	11/30/2015	1.29 \$M		FALSE	2	2	1165 sqft		N/A
15	15063340	44444444	37799474	-122414835	11/30/2015	1.29 \$M		FALSE	2	2	1165 sqft		N/A
16	15063340	44444444	37799474	-122414835	11/30/2015	1.29 \$M		FALSE	2	2	1165 sqft		N/A

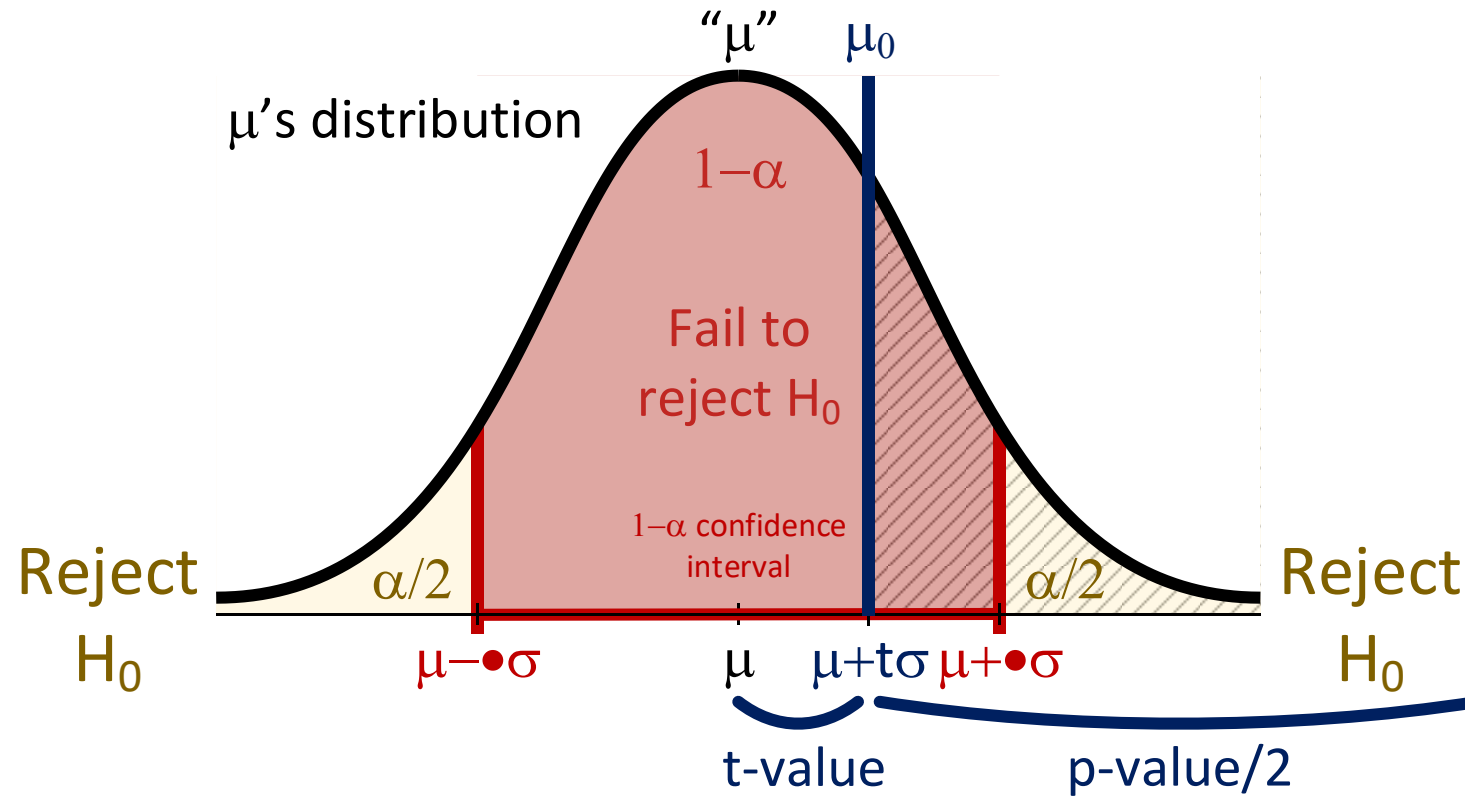
⑤ REFINE the Data



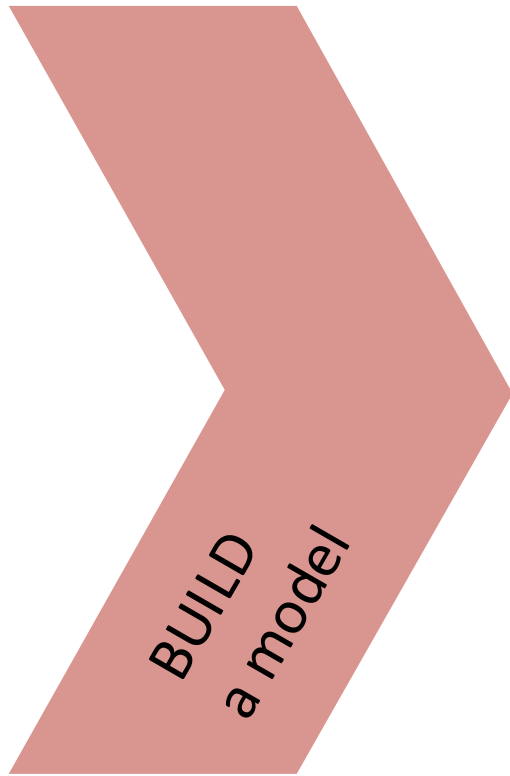
- Refine the Data
 - Identify trends and outliers
 - Apply descriptive and inferential statistics
 - Document and transform data

5 REFINE the Data

We will apply inferential statistics



⑥ BUILD a Model



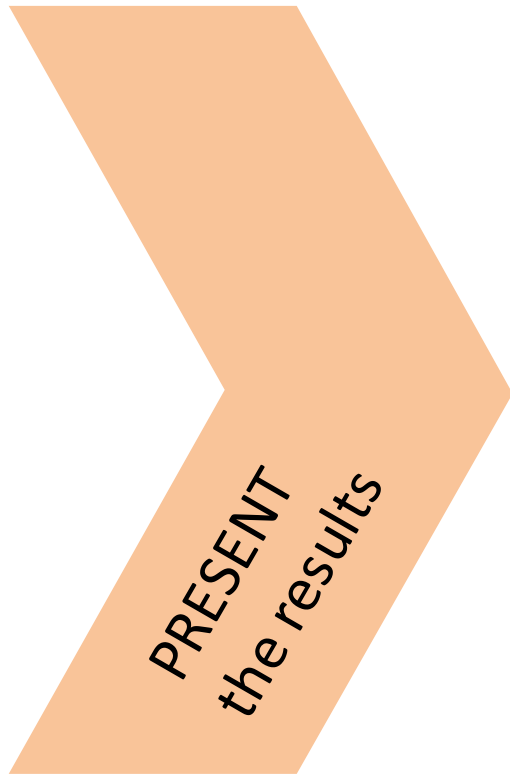
- Build a Model
 - Select appropriate model
 - Build model
 - Evaluate and refine model

⑥ BUILD a Model

Types of machine learning algorithms we will study in this course

	Continuous	Categorical
Supervised (a.k.a., predictive modeling)	k-Nearest Neighbors <i>(session 5)</i> Linear Regression <i>(sessions 8 and 9)</i> Trees <i>(session 16)</i> Time Series <i>(session 19)</i>	k-Nearest Neighbors <i>(session 5)</i> Logistic Regression <i>(session 11)</i> Trees <i>(session 16)</i> Natural Language Processing <i>(session 18)</i>
Unsupervised	Clustering <i>(session 14)</i> Natural Language Processing <i>(session 18)</i>	

⑦ 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

7 PRESENT the Results

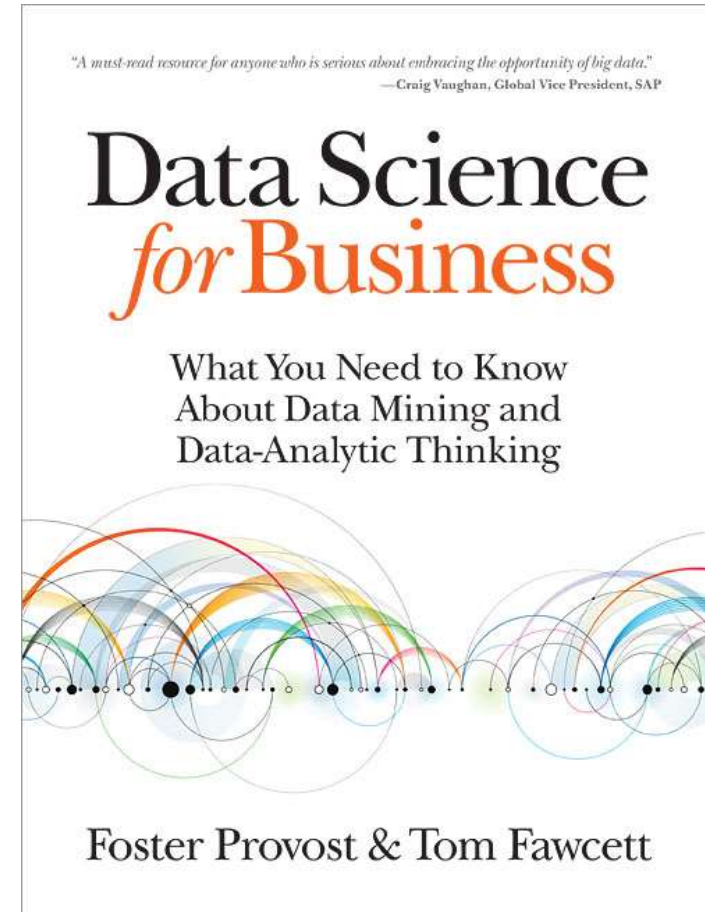
Know Your Audience



Corina Rosu © 123RF.com

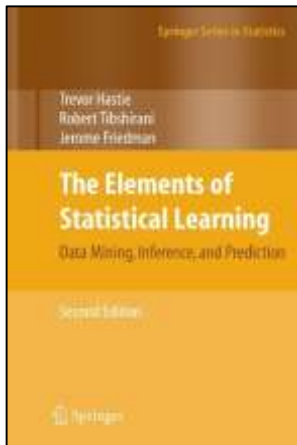
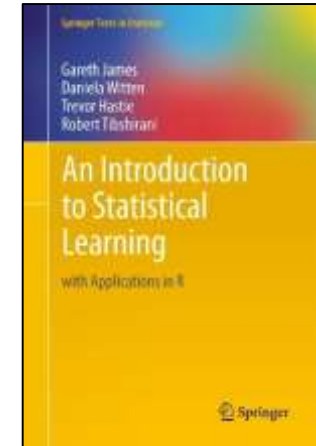
A great resource to follow along the class (or afterwards...) (*will mostly reference pre-class reading materials; optional; not required but recommended for the course*)

- Data Science for Business (by Provost and Fawcett) [[Link](#)]
 - General Assembly holds several copies in its library



A couple of resources to follow along the class (or afterwards...) (*will mostly reference post-class reading materials; 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](#)



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

DS

① IDENTIFY the Problem

1 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
- The SMART Goals Framework

By asking a good question and setting a clear aim:



- You set yourself up for success
 - “A problem well stated is half solved” – Charles Kettering
- You help other data scientists learn from and reproduce your work
 - You establish the basis for making your analysis reproducible
- You also help them expand on your work in the future

The SMART Goals Framework for Data Science

(https://en.wikipedia.org/wiki/SMART_criteria)

S _{PECIFIC}	The dataset and key variables are clearly defined
M _{EASURABLE}	The type of analysis and major assumptions are articulated
A _{TTAINABLE}	The question you are asking is feasible for your dataset and is not likely to be biased
R _{EPRODUCIBLE}	Another person (or you in 6 months!) can read your state and understand exactly how your analysis is performed
T _{IME-BOUND}	You clearly state the time period and population for which this analysis will pertain

Trends often change over time and vary by the population of source of your data. It is important to clearly define who/what you included in your analysis as well as the time period for the analysis

DS

7 PRESENT the Results

We built a model! Now what?

- We've built our model, but there is still a gap between our Jupyter Notebook with its plots and figures and a slideshow needed to present our results
- The course focus on two core concepts:
 - Developing consistent practices
 - Interpreting metrics to evaluate and improve model performance
- But what does that mean to your audience?

We built a model! Now what? (cont.)

- Imagine how a non-technical audience might respond to the following statements:
 - “The predictive model I built has an accuracy of 80%”
 - “Logistic regression was optimized with L2 regularization”
 - “Gender was more important than age in the predictive model because it has a ‘larger coefficient’”
 - “Here’s the AUC chart that shows how well the model did”

We built a model! Now what? (cont.)

- Who is your audience? Are they technical? What are their concerns?
 - In a business setting, you may be the only person who can interpret what you've built
- Some people may be familiar with basic visualization, but you will likely have to do a lot of “hand holding”
- You need to be able to efficiently explain your results in a way that makes sense to all stakeholders (technical or not)

We built a model! Now what? (cont.)

- In this section, we'll focus on communicating results for “simpler” problems, but this applies to any type of model you may work with

Showing our Work

- We've spent a lot of time exploring our data and building a reasonable model that performs well
- However, if we look at our visuals, they are most likely:

- Statistically heavy:
most people don't
understand histograms

- Overly complicated:
scatter matrices
produce too much
information

- Poorly labeled: code
doesn't require adding
labels, so you may not
have added them

To convey important information to your audience, make sure your charts are simplified, easily interpretable, and clearly labeled

Simplified

- At most, you'll want to include figures that either explain a variable on its own or explain that variable's relationship against a target
- If your model used a data transformation (like natural log), just visualize the original data
- Remove unnecessary complexity

Easily interpretable

- Any stakeholder looking at a figure should be seeing the exact same thing you're seeing
 - A good test for this is to share the visual with others less familiar with the data and see if they come to the same conclusion
 - How long did it take them?

Clearly labeled

- Take the time to clearly label your axis, title your plot, and double check your scales – especially if the figures should be comparable
- If you're showing two graphs side by side, they should follow the same Y axis

When building visuals for another audience, ask yourself who, what, and how

Who

- Who is my target audience for the visual?

What

- What do they already know about this project?
What do they need to know?

How

- How does my project affect this audience? How might they interpret (or misinterpret) the data?

Visualizing Models over Variables

- One effective way to explain your model over particular variables is to plot the predicted values against the most explanatory variables
- E.g., in logistic regression, plotting the probability of a class against a variable can help explain the range of effect of the model

Visualizing Performance against Baseline

- Another approach of visualization is the effect of your model against a baseline, or – even better – against previous models
- Plots like this will also be useful when talking to your peers – other data scientists or analysts who are familiar with your project and interested in the progress you've made

Slides © 2016 Ivan Corneillet Where Applicable
Do Not Reproduce Without Permission