# GENERAL ASSEMBLY

# Welcome to General Assembly

## *Part-Time Data Science*

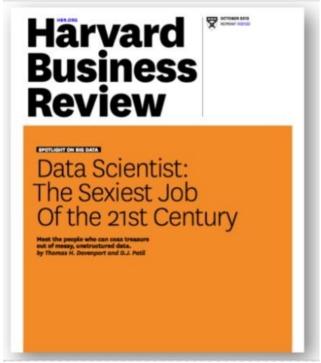**Schedule**    7:00 PM: Introductory Slides
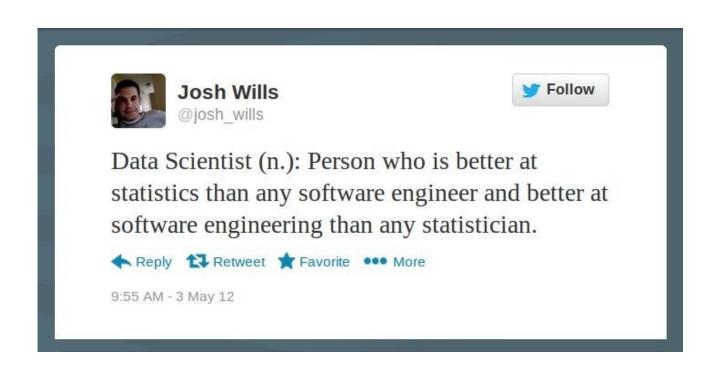8:00 PM: Trello, Github,
Anaconda
8:30 PM: Pivot to Python

# WHAT IS DATA SCIENCE?

# WHAT IS DATA SCIENCE?



Job Trends from Indeed.com
— big-data — data-science



Harvard Business Review
OCTOBER 2012

SPOTLIGHT ON BIG DATA

Data Scientist:
The Sexiest Job
Of the 21st Century

Meet the people who can coax treasure
out of messy, unstructured data.
by Thomas H. Davenport and D.J. Patil

# WHAT IS DATA SCIENCE?

**Josh Wills**
@josh_wills

Follow

Data Scientist (n.): Person who is better at statistics than any software engineer and better at software engineering than any statistician.

Reply   Retweet   Favorite   More

9:55 AM - 3 May 12

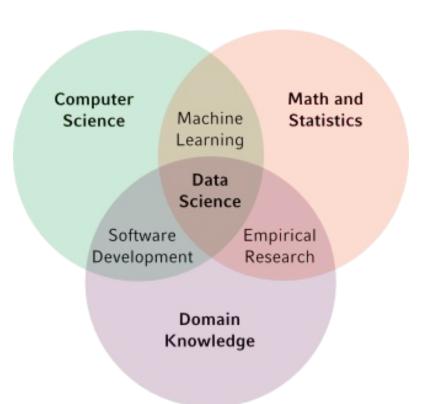# WHAT IS DATA SCIENCE?

# BY FIELD

# WHAT IS DATA SCIENCE?

- **How They're Using Data Science:**

- Airbnb prioritizes listings in popular areas, making desirable Airbnbs easier for users to find.
- KAID Health uses natural language processing to mine clinical notes, allowing providers to find patients for clinical trials
- UPS optimizes package drop-off and delivery transport using machine learning and AI to predict delivery obstacles (e.g., weather, traffic).

# WHAT IS DATA SCIENCE?

**Consider these products and services:**

- How do they utilize data science?
- What kinds of data do you think they use?
- How might they leverage data science in other parts of their business?

# DATA SCIENCE WORKFLOW

1. <u>Identify</u> the problem

2. <u>Acquire</u> the data

3. <u>Parse</u> the data

4. <u>Mine</u> the data

5. <u>Refine</u> the data

6. <u>Build</u> a data model

7. <u>Present</u> the results

**DATA SCIENCE WORKFLOW**

**IDENTIFY THE PROBLEM**
- ☐ Identify business/product objectives
- ☐ Identify and hypothesize goals and criteria for success
- ☐ Create a set of questions for identifying correct data set

**ACQUIRE THE DATA**
- ☐ Identify the "right" data set(s)
- ☐ Import data and set up local or remote data structure
- ☐ Determine most appropriate tools to work with data

**PARSE THE DATA**
- ☐ Read any documentation provided with the data
- ☐ Perform exploratory data analysis
- ☐ Verify the quality of 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)

**REFINE THE DATA**
- ☐ Identify trends and outliers
- ☐ Apply descriptive and inferential statistics
- ☐ Document and transform data

**BUILD A DATA MODEL**
- ☐ Select appropriate model
- ☐ Build model
- ☐ Evaluate and refine model

**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

*Identify · Acquire · Parse · Mine · Refine · Build · Present*

# IDENTIFY

‣ **Why are you doing this in the first place?**

‣ *We believe there is a market for automating detailed medical forecasts for individual claims*

‣ **Who are the stakeholders?**

‣ **What data will you need?**

‣ *Is it available to us?*

‣ *Is is public or proprietary? Is your work easily duplicated if the former?*

‣ **How will you define success?**



**DATA SCIENCE WORKFLOW**

**IDENTIFY THE PROBLEM**
☐ Identify business/product objectives
☐ Identify and hypothesize goals and criteria for success
☐ Create a set of questions for identifying correct data set

**ACQUIRE THE DATA**
☐ Identify the "right" data set(s)
☐ Import data and set up local or remote data structure
☐ Determine most appropriate tools to work with data

**PARSE THE DATA**
☐ Read any documentation provided with the data
☐ Perform exploratory data analysis
☐ Verify the quality of 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)

**REFINE THE DATA**
☐ Identify trends and outliers
☐ Apply descriptive and inferential statistics
☐ Document and transform data

**BUILD A DATA MODEL**
☐ Select appropriate model
☐ Build model
☐ Evaluate and refine model

**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

# WHAT IS DATA SCIENCE?

# ACQUIRE

‣ **Can you supplement your data?**

‣ *Is there information in the clinical notes that might not appear in bills?*

‣ **How is it stored?**

‣ *CBI stores data in relational databases and .csv files*

‣ **What tools will you need to work with it?**

‣ *Software (to manipulate data; fit algorithms) and hardware (to handle computations)*

**DATA SCIENCE WORKFLOW**

**IDENTIFY THE PROBLEM**
- ☐ Identify business/product objectives
- ☐ Identify and hypothesize goals and criteria for success
- ☐ Create a set of questions for identifying correct data set

*Identify*

**ACQUIRE THE DATA**
- ☐ Identify the "right" data set(s)
- ☐ Import data and set up local or remote data structure
- ☐ Determine most appropriate tools to work with data

*Acquire*

**PARSE THE DATA**
- ☐ Read any documentation provided with the data
- ☐ Perform exploratory data analysis
- ☐ Verify the quality of the data

*Parse*

**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)

*Mine*

**REFINE THE DATA**
- ☐ Identify trends and outliers
- ☐ Apply descriptive and inferential statistics
- ☐ Document and transform data

*Refine*

**BUILD A DATA MODEL**
- ☐ Select appropriate model
- ☐ Build model
- ☐ Evaluate and refine model

*Build*

**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

*Present*

# PARSE

‣ **How do you get raw data into a format you can work with?**

‣ *This is the purview of "Data Engineers"*

‣ **What documentation is available, if any?**

‣ *Data dictionaries are ideal*

‣ **How much munging will it require?**

‣ *Are all your date fields valid dates?*

‣ *Are some fields mysteriously empty?*



**DATA SCIENCE WORKFLOW**

**IDENTIFY THE PROBLEM**
☐ Identify business/product objectives
☐ Identify and hypothesize goals and criteria for success
☐ Create a set of questions for identifying correct data set

**ACQUIRE THE DATA**
☐ Identify the "right" data set(s)
☐ Import data and set up local or remote data structure
☐ Determine most appropriate tools to work with data

**PARSE THE DATA**
☐ Read any documentation provided with the data
☐ Perform exploratory data analysis
☐ Verify the quality of 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)

**REFINE THE DATA**
☐ Identify trends and outliers
☐ Apply descriptive and inferential statistics
☐ Document and transform data

**BUILD A DATA MODEL**
☐ Select appropriate model
☐ Build model
☐ Evaluate and refine model

**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

*Identify* *Acquire* *Parse* *Mine* *Refine* *Build* *Present*

# MINE & REFINE

‣ **Combining and Transforming the data**

‣ *Aggregating inpatient and outpatient bills to yield a longitudinal service history for individual claimants*

‣ **Mining the data to find predictive insights**

‣ *Example: what's a predictor of shoulder surgery?*

| Claimant | Age | Torn Rotator Cuff | Shoulder Sprain | Rotator Cuff Surgery |
|----------|-----|-------------------|-----------------|----------------------|
| Jim | 28 | True | True | True |
| Pat | 40 | False | True | False |

**DATA SCIENCE WORKFLOW**

**IDENTIFY THE PROBLEM**
- ☐ Identify business/product objectives
- ☐ Identify and hypothesize goals and criteria for success
- ☐ Create a set of questions for identifying correct data set

*Identify*

**ACQUIRE THE DATA**
- ☐ Identify the "right" data set(s)
- ☐ Import data and set up local or remote data structure
- ☐ Determine most appropriate tools to work with data

*Acquire*

**PARSE THE DATA**
- ☐ Read any documentation provided with the data
- ☐ Perform exploratory data analysis
- ☐ Verify the quality of the data

*Parse*

**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)

*Mine*

**REFINE THE DATA**
- ☐ Identify trends and outliers
- ☐ Apply descriptive and inferential statistics
- ☐ Document and transform data

*Refine*

**BUILD A DATA MODEL**
- ☐ Select appropriate model
- ☐ Build model
- ☐ Evaluate and refine model

*Build*

**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

*Present*

# WHAT IS DATA SCIENCE?

# MODEL

‣ What model or models are most appropriate for the data and the problem?

‣ *Is your data linear (each additional square foot yields a higher price) or non-linear (small houses in desirable neighborhoods cost more than large houses in undesirable ones)*

‣ How can you be sure your model results generalize?

‣ *We need to evaluate out-of-sample data that our model(s) haven't trained on to understand this*



**DATA SCIENCE WORKFLOW**

**IDENTIFY THE PROBLEM**
☐ Identify business/product objectives
☐ Identify and hypothesize goals and criteria for success
☐ Create a set of questions for identifying correct data set

**ACQUIRE THE DATA**
☐ Identify the "right" data set(s)
☐ Import data and set up local or remote data structure
☐ Determine most appropriate tools to work with data

**PARSE THE DATA**
☐ Read any documentation provided with the data
☐ Perform exploratory data analysis
☐ Verify the quality of 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)

**REFINE THE DATA**
☐ Identify trends and outliers
☐ Apply descriptive and inferential statistics
☐ Document and transform data

**BUILD A DATA MODEL**
☐ Select appropriate model
☐ Build model
☐ Evaluate and refine model

**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

# PRESENT

‣ **What narrative do I want to tell?**

‣ *Why does the model predict claimant x will or won't get surgery?*

‣ **What inherent limitations should be disclosed?**

‣ *Did you collect the data your models are based on...*

‣ *Will your model need to be retrained?*

‣ **Were the criteria for success met?**

‣ *With more time, how would we improve our result?*



**DATA SCIENCE WORKFLOW**

**IDENTIFY THE PROBLEM**
- ☐ Identify business/product objectives
- ☐ Identify and hypothesize goals and criteria for success
- ☐ Create a set of questions for identifying correct data set

**ACQUIRE THE DATA**
- ☐ Identify the "right" data set(s)
- ☐ Import data and set up local or remote data structure
- ☐ Determine most appropriate tools to work with data

**PARSE THE DATA**
- ☐ Read any documentation provided with the data
- ☐ Perform exploratory data analysis
- ☐ Verify the quality of 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)

**REFINE THE DATA**
- ☐ Identify trends and outliers
- ☐ Apply descriptive and inferential statistics
- ☐ Document and transform data

**BUILD A DATA MODEL**
- ☐ Select appropriate model
- ☐ Build model
- ☐ Evaluate and refine model

**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

# WHAT IS DATA SCIENCE?

# (DEPLOY)

‣ **How do we run the model in production?**

‣ *I pickle the model*

‣ *When a user calls our API, the model is loaded*

‣ *But we have to provide the same features to our model as I did fitting it, so the data submitted on a web form has to go through a series of transformations first. (For example, taking a diagnosis and passing a rate into the model.)*

**DATA SCIENCE WORKFLOW**

**IDENTIFY THE PROBLEM**
☐ Identify business/product objectives
☐ Identify and hypothesize goals and criteria for success
☐ Create a set of questions for identifying correct data set

*Identify*

**ACQUIRE THE DATA**
☐ Identify the "right" data set(s)
☐ Import data and set up local or remote data structure
☐ Determine most appropriate tools to work with data

*Acquire*

**PARSE THE DATA**
☐ Read any documentation provided with the data
☐ Perform exploratory data analysis
☐ Verify the quality of the data

*Parse*

**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)

*Mine*

**REFINE THE DATA**
☐ Identify trends and outliers
☐ Apply descriptive and inferential statistics
☐ Document and transform data

*Refine*

**BUILD A DATA MODEL**
☐ Select appropriate model
☐ Build model
☐ Evaluate and refine model

*Build*

**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

*Present*

# MODELING

# MODELING

# Video #1 - Margaritas

‣ https://www.youtube.com/watch?v=t_3fnVqNOUc

# Video #2 - Tennis

‣ https://www.youtube.com/watch?v=eKD5gxPPeY0

## Comparing these examples

- **How did we cluster the margaritas together?**
- **How did we make a prediction about tennis?**
- **What seems different about the approaches?**

# Features & Target Variables

‣ $x$: features (inputs)

‣ $y$: target variable (output; can be numeric or binary)

‣ Models: values in, value(s) out

$$y = f(x_1, x_2, x_3...)$$

# TARGET VARIABLE FLAVORS

| continuous | categorical |
| --- | --- |
| Height of children | Eye colors |
| Weight of cars | Courses at GA |
| Speed of the train | Highest degree |
| Temperature | Gender |
| Stock price | If an email is spam or not |

# MACHINE LEARNING PROBLEMS

|  | continuous | categorical |
|---|---|---|
| **supervised** | regression | classification |
| **unsupervised** | dimension reduction | clustering |

# REGRESSION
## (CONTINUOUS, SUPERVISED)

‣ Build a model to predict a continuous value that best fits data

‣ Minimize error without overfitting

‣ *Example: Linear Regression*

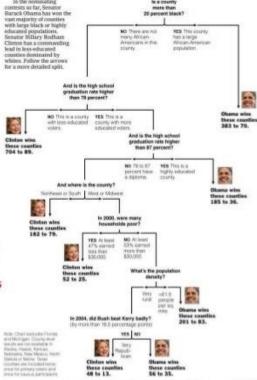http://setosa.io/ev/ordinary-least-squares -regression/

# CLASSIFICATION
## (CATEGORICAL, SUPERVISED)

‣ Map features to categorical target classes.



Decision Tree: The Obama-Clinton Divide

Can we learn how counties vote?

New York Times
April 16, 2008

Decision Trees:
a sequence of tests.
Representation very natural for humans.
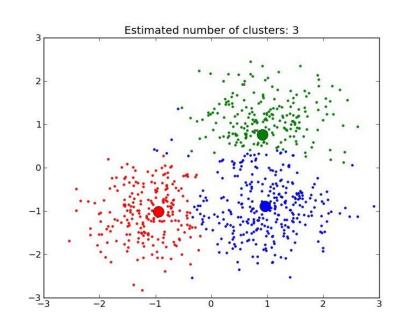Style of many "How to" manuals and trouble-shooting procedures.

# CLUSTERING
## (CATEGORIAL, UNSUPERVISED)

‣ Purpose is representation – anything that helps you better understand the data

‣ Finding common threads that a human couldn't see (imagine 10 or 100 or 1000 inputs)

‣ *Example: K-Means*

https://www.naftaliharris.com/blog/visualizing-k-means-clustering/

Estimated number of clusters: 3

# Regression: How Much Physical Therapy?

| Claimant | Split | Severe Underlying Injury | Surgery After Injury | PT Units |
|----------|-------|--------------------------|----------------------|----------|
| Jim | Train | Yes | Yes | 13 |
| Pat | Train | No | No | 4 |
| Dimitri | Test | Yes | No | **?** |

# Classification: Likelihood of Major Surgery

| Claimant | Split | Rate That Diagnosis Receives Surgery | Proximity to Accident (years) | Surgical Likelihood |
|----------|-------|--------------------------------------|-------------------------------|---------------------|
| Deb | Training | .5 | <1 | 100% |
| Jack | Training | .2 | 3 | 0% |
| Dimitri | Test | .15 | <1 | **?** |

# CASE STUDY

| Item | Orders Since Customer Included Item | Customer Order Rate | Order Probability |
|------|-----------------------------------|---------------------|-------------------|
| Bananas | 1 | 90% | ? |
| Baking Soda | 6 | 5% | ? |

## With (a) partner(s), map out where this fits in the data science workflow – high level:
what phase is it, what comes before it, what comes after it

HOWEVER

There is a lot of variation in what data scientists do. Most of my experience is… *not* modeling.

- I write rules to e.g. deidentify documents
- I label data
- I've built dashboards
- I've handled ETL
- I've migrated data
- I've contributed to research

# GO FURTHER

## Understanding Exploratory Data Analysis (EDA)

https://www.kaggle.com/pmarcelino/comprehensive-data-exploration-with-python

## Annotated Titanic Workflow

https://www.kaggle.com/headsortails/pytanic

# That's a wrap!