

WELCOME TO DATA SCIENCE

Trevor Lindsay

Data Scientist, Facebook

WELCOME TO DATA SCIENCE

LEARNING OBJECTIVES

- Describe the roles and components of a successful learning environment
- Define data science and the data science workflow
- Apply the data science workflow to meet your classmates
- Setup your development environment and review python basics

DATA SCIENCE

WELCOME TO GA!

WELCOME TO GA!

- General Assembly is a global community of individuals empowered to pursue the work we love.
- General Assembly's mission is to build our community by transforming millions of thinkers into creators.

YOUR INSTRUCTOR

Trevor Lindsay

Lead Instructor
Data Scientist, Facebook
trevor.lindsay@yahoo.com



STUDENT SERVICES



Matt Jones

Student Experience Associate

studentservicesSF@ga.co

Course Logistics + Campus Questions

OTHERS YOU MAY SEE



RAY HSIA

Instructor Manager



NIÑA PINEDA

Front Lines Lead



VANESSA OHTA

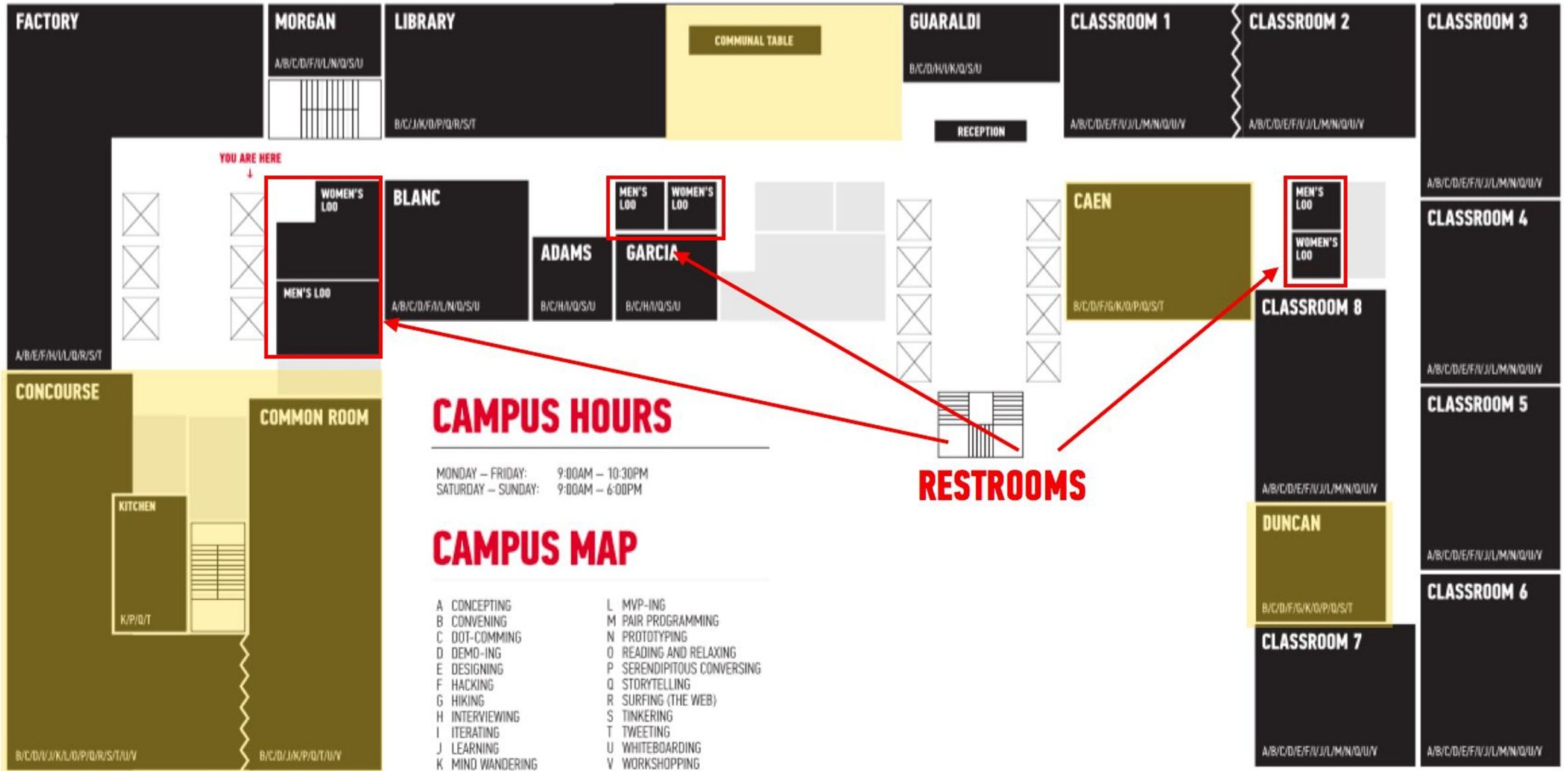
Senior
Instruction Manager

LET'S GET TO KNOW EACH OTHER

- Introduce yourself (name, work experience, where in the Bay Area you live, etc.)
- Tell us why you are taking this course and what you hope to get out of it
- Open your preferred music app (Spotify, Apple Music, etc.) and share the last song you listened to

COME WORK ON CAMPUS!

- Hours
 - 8am - 10pm, Monday to Friday
 - 9am - 6pm, Saturday and Sunday
- Reception can help with:
 - Loaner equipment
 - Lost and found
 - Free coffee and tea



PUBLIC USE SPACES

DATA SCIENCE

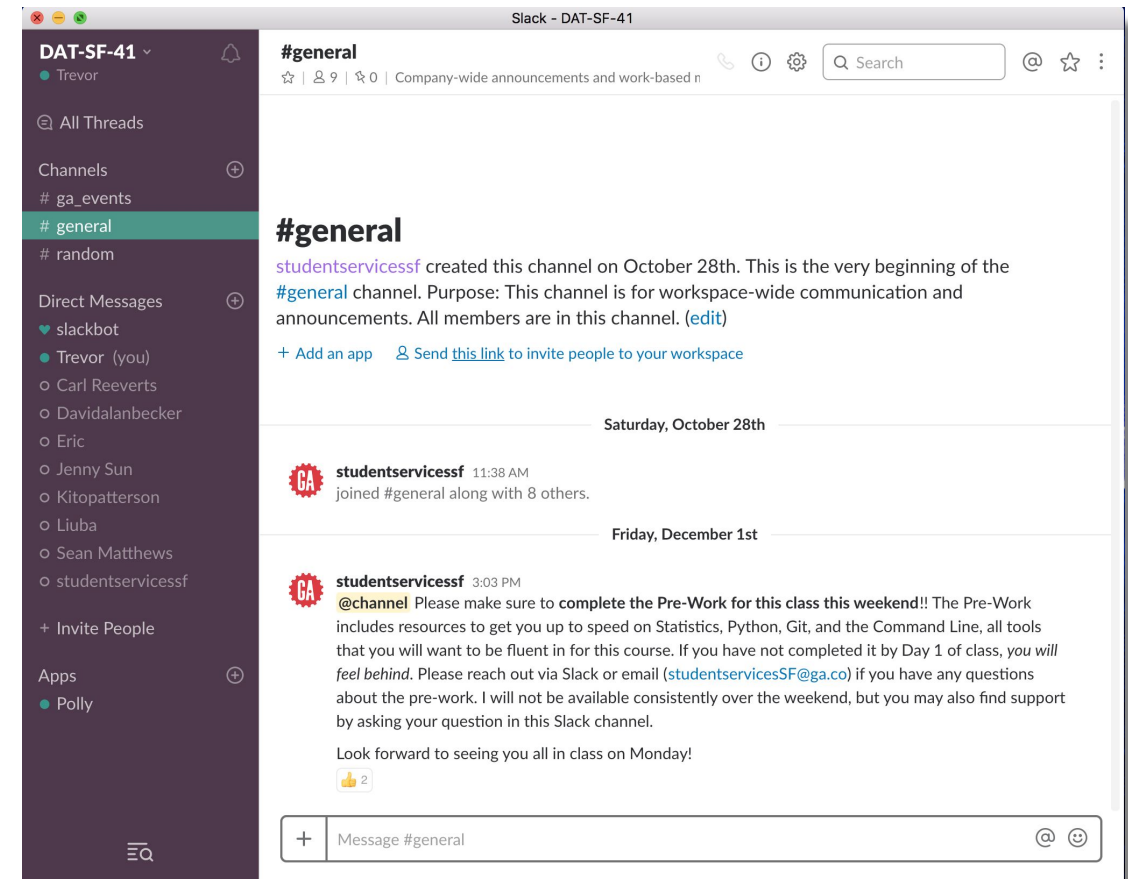
Student Experience

SLACK



dat-sf-41.slack.com

All course communication with each other and instructors will happen here

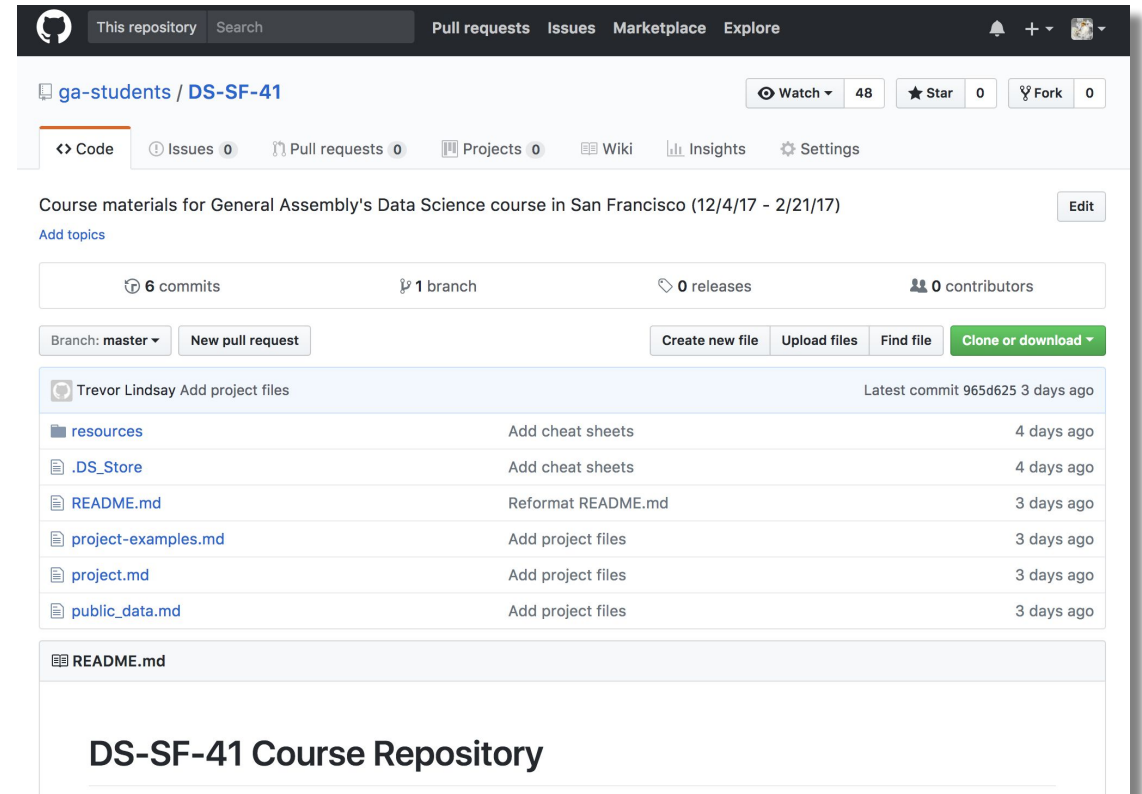


GITHUB

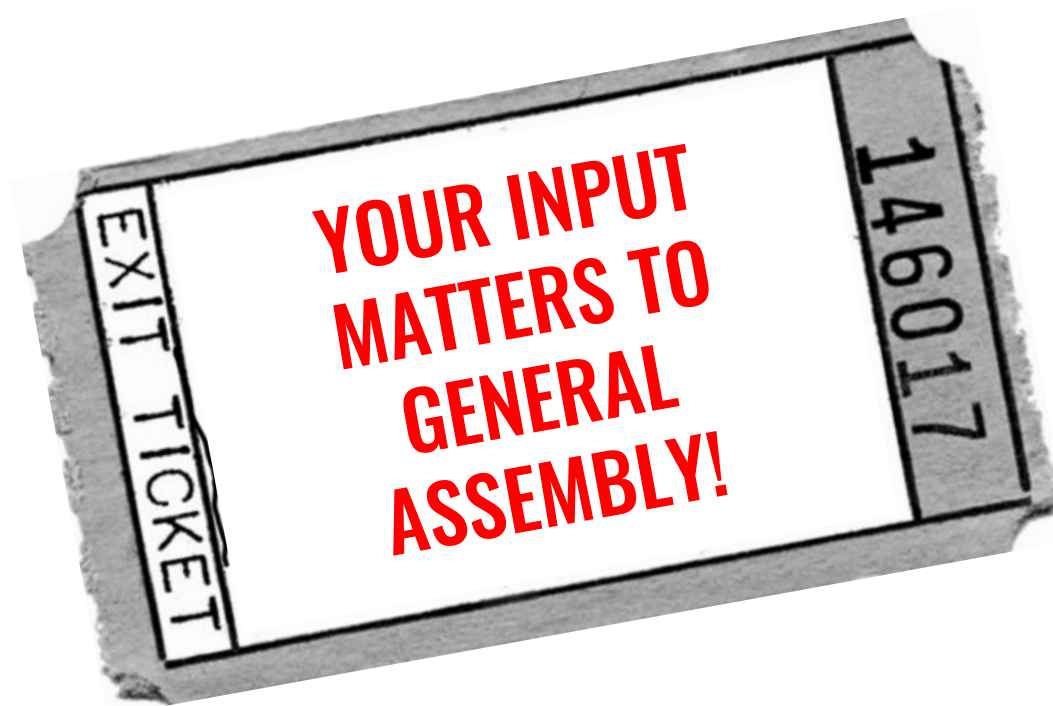


github.com/ga-students/DS-SF-41

All the course resources (sample code, assignments, slides) you need will be stored here.



EXIT TICKETS



<https://tinyurl.com/dat-sf-41>

At the end of each class, please fill out an exit ticket to provide feedback on your experience

CLASSROOM CULTURE



Let's all agree to:

- Treat each other with respect
- Avoid distractions during class
- Actively participate
- Arrive on time

INSTRUCTOR PHILOSOPHY

- Seek an optimal pace
- Facilitate healthy and active learning environment
- Involve everyone
- Communicate early and often
 - Feedback should be a busy two-way street

CONTENT PHILOSOPHY

- Alignment with industry
- Understand foundational principles of theories and ideas
- Balance depth with breadth
- Engaging and interesting data
 - Examples of real world data science
- Course project

KEYS TO SUCCESS

- Effort > prior knowledge
- Ask questions **early** and **often**
- Keep track of what you've learned and overall progress
- Use each other
- Be patient
- Time management

CLASS STRUCTURE

- Lecture + Theory: Introduction to material
 - 45 min - 1 hr
- Interactive Coding: Demonstrate material
 - 1 hr
- Class Work: Practice material on real data
 - 30 min - 1 hr
- 5-10 min breaks to digest the material and reset

CLASS STRUCTURE (CONT.)

- All of the course material is in the Github repo
- Almost all of the classes take place in Jupyter Notebooks
- Office hours from 5:30pm - 6:30pm
- Assignments (required for graduation)
 - 2 homeworks
 - 1 final project
 - 4 small unit projects (optional)
- Use Slack to communicate with classmates and me
- If you're late or have to miss class, please inform me at your earliest convenience

GA GRADUATION REQUIREMENTS

HOMEWORK
(COMPLETE 80% OF
HOMEWORK/LABS)

ATTENDANCE
(MISS NO MORE THAN 2
CLASSES)

**FINAL
PROJECT**

**COMMUNITY
ENGAGEMENT**
PARTICIPATION +
FEEDBACK

DATA SCIENCE

PRE-WORK

PRE-WORK REVIEW

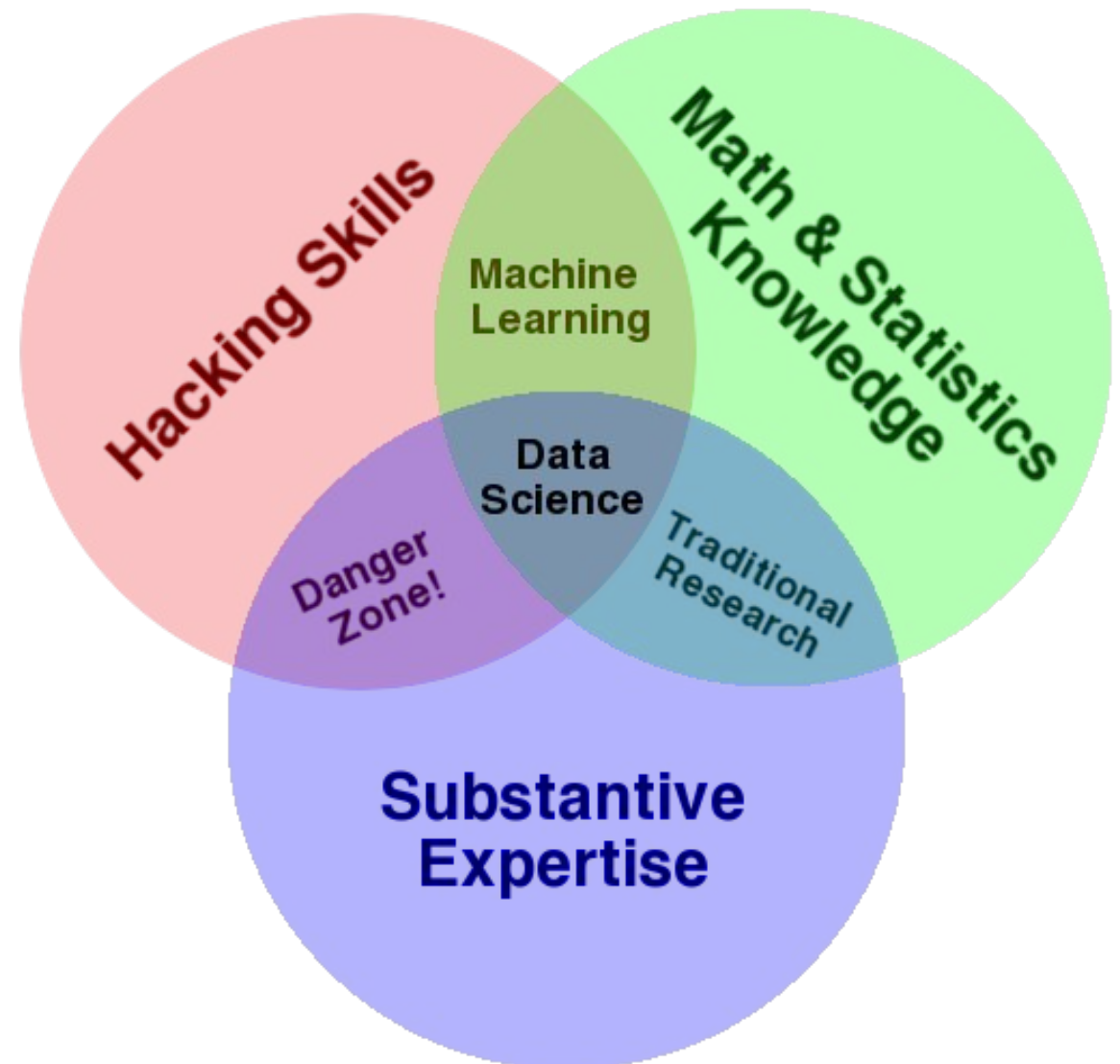
- Define basic data types used in object-oriented programming
- Recall the Python syntax for lists, dictionaries, and functions
- Create files and navigate directories using the command line interface

INTRODUCTION

WHAT IS DATA SCIENCE?

WHAT IS DATA SCIENCE?

- A set of tools and techniques for data
- Interdisciplinary problem-solving
- Application of scientific techniques to practical problems



WHO IS A DATA SCIENTIST?

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

When Jonathan Goldman arrived for work in June 2006 at LinkedIn, the business networking site, the place still felt like a start-up. The company had just under 8 million accounts, and the number was growing quickly as existing members invited their friends and colleagues to join. But users weren't seeking out connections with the people who were already on the site at the rate executives had expected. Something was apparently missing in the social experience. As one LinkedIn manager put it, "It was like arriving at a conference reception and realizing you don't know anyone. So you just stand in the corner sipping your drink—and you probably leave early."

70 Harvard Business Review October 2012

MODERN DATA SCIENTIST

Data Scientist, the sexiest job of the 21st century, requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

MATH & STATISTICS

- ☆ Machine learning
- ☆ Statistical modeling
- ☆ Experiment design
- ☆ Bayesian inference
- ☆ Supervised learning: decision trees, random forests, logistic regression
- ☆ Unsupervised learning: clustering, dimensionality reduction
- ☆ Optimization: gradient descent and variants



PROGRAMMING & DATABASE

- ☆ Computer science fundamentals
- ☆ Scripting language e.g. Python
- ☆ Statistical computing packages, e.g., R
- ☆ Databases: SQL and NoSQL
- ☆ Relational algebra
- ☆ Parallel databases and parallel query processing
- ☆ MapReduce concepts
- ☆ Hadoop and Hive/Pig
- ☆ Custom reducers
- ☆ Experience with xaaS like AWS

DOMAIN KNOWLEDGE & SOFT SKILLS

- ☆ Passionate about the business
- ☆ Curious about data
- ☆ Influence without authority
- ☆ Hacker mindset
- ☆ Problem solver
- ☆ Strategic, proactive, creative, innovative and collaborative

COMMUNICATION & VISUALIZATION

- ☆ Able to engage with senior management
- ☆ Story telling skills
- ☆ Translate data-driven insights into decisions and actions
- ☆ Visual art design
- ☆ R packages like ggplot or lattice
- ☆ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau

MarketingDistillery.com is a group of practitioners in the area of e-commerce marketing. Our fields of expertise include: marketing strategy and optimization; customer tracking and on-site analytics; predictive analytics and econometrics; data warehousing and big data systems; marketing channel insights in Paid Search, SEO, Social, CRM and brand.

Marketing
DISTILLERY
(c) Krzysztof Zawadzki

WHO IS A DATA SCIENTIST?

“...unique blend of skills that can both unlock the insights of data and tell a fantastic story via the data.”

- *DJ Patil, Former U.S. Chief Data Scientist*

“...someone who blends, math, algorithms, and an understanding of human behavior with the ability to hack systems together to get answers to interesting human questions from data.”

- *Hilary Mason, Data Scientist in Residence at Accel*

“...better at statistics than any software engineer and better at software engineering than any statistician.”

- *Unknown*

DATA SCIENCE IS EVERYWHERE

NETFLIX



FiveThirtyEight



amazon.com[®]



Google

Walmart 
Save money. Live better.

FACEBOOK DATA SCIENTIST

- Apply your expertise in quantitative analysis, data mining, and the presentation of data to see beyond the numbers and understand how our users interact with both our consumer and business products
- Partner with Product and Engineering teams to solve problems and identify trends and opportunities
- Inform, influence, support, and execute our product decisions and product launches
- Work across four main areas:
 - Product Operations
 - Exploratory Data Analysis
 - Product Leadership
 - Data Infrastructure

GOOGLE DATA SCIENTIST

- Evaluate and improve Google's products by collaborating with a multi-disciplinary team of engineers and analysts on a wide range of problems
- Bring analytical rigor and statistical methods to the challenges of measuring quality, improving consumer products, and understanding the behavior of end-users, advertisers, and publishers
- Work with large, complex data sets. Solve difficult, non-routine analysis problems, applying advanced analytical methods as needed.
- Conduct end-to-end analysis that includes data gathering and requirements specification, processing, analysis, ongoing deliverables, and presentations.

UBER DATA SCIENTIST

- Familiarity with technical tools for analysis - Python (with Pandas, etc.), R, SQL
- Programming chops - demonstrable familiarity (work experience, Github account) with programming concepts. Python skills and previous software engineering background a plus
- Research mindst - ability to structure a project from idea to experimentation to prototype to implementation
- Driven and focused self-starters, great communicators, amazing follow-through - you passionately pursue your work and love the responsibility of being individually empowered
- A preference for quality over quantity - you get the math right and aspire to build the right solution; you like a team that holds each other to a high bar

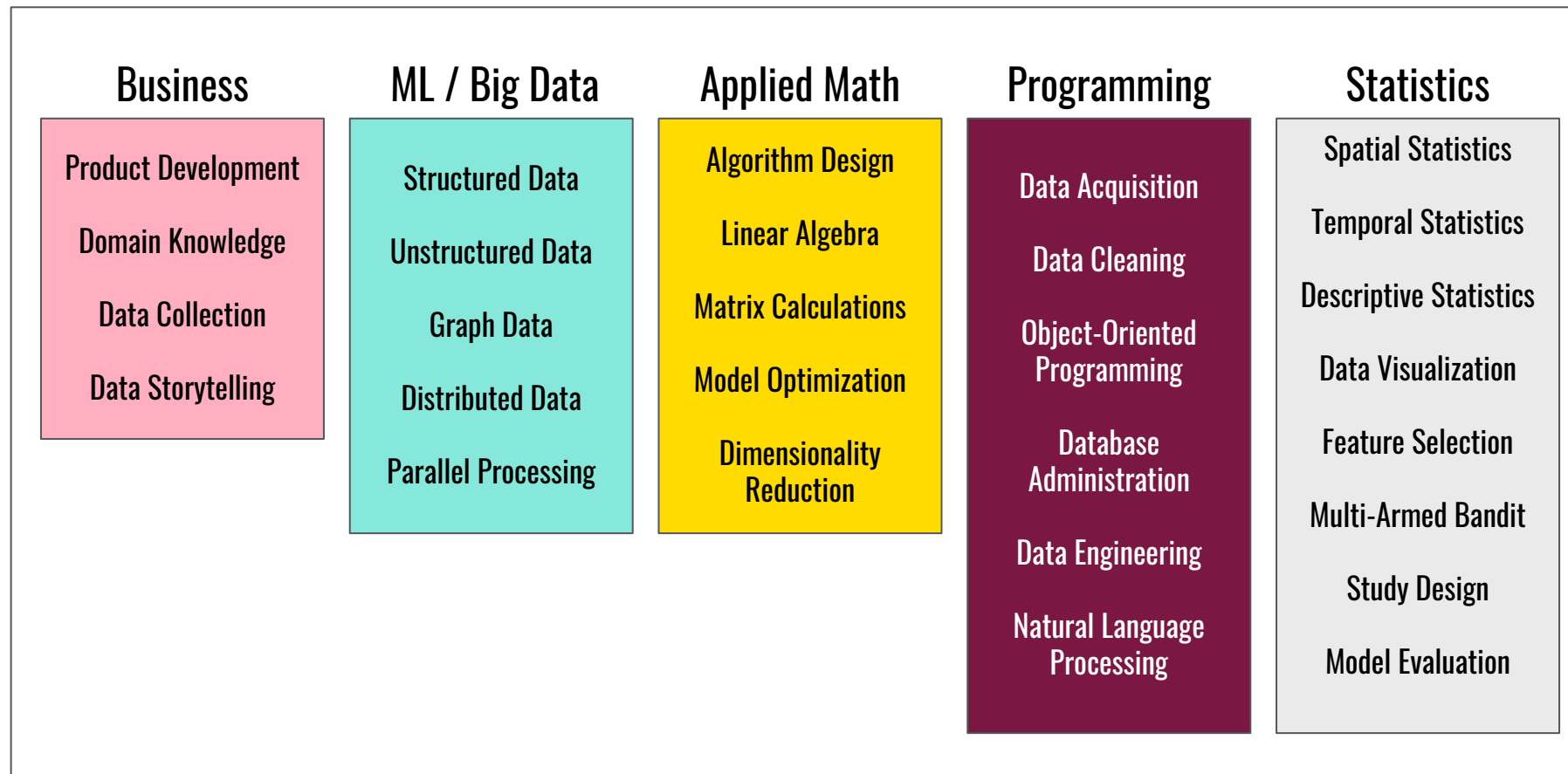
WHAT ARE THE ROLES IN DATA SCIENCE?

- Data Science involves a variety of roles, not just one.

Data Developer	Developer	Engineer	
Data Researcher	Researcher	Scientist	Statistician
Data Creative	Jack of All Trades	Artist	Hacker
Data Businessperson	Leader	Businessperson	Entrepreneur

WHAT ARE THE ROLES IN DATA SCIENCE?

- Data Science involves a variety of skill sets, not just one.



QUIZ

DATA SCIENCE BASELINE

ACTIVITY: DATA SCIENCE BASELINE QUIZ



EXERCISE

DIRECTIONS (10 minutes)

1. Form groups of three.
2. Answer the following questions.
 - a. True or False: Gender (coded male=0, female=1) is a continuous variable.
 - b. According to the table on the next slide, BMI is the _____
 - i. Outcome
 - ii. Predictor
 - iii. Covariate
 - c. Draw a normal distribution
 - d. True or False: Linear regression is an unsupervised learning algorithm.
 - e. What is a hypothesis test?

ACTIVITY: DATA SCIENCE BASELINE QUIZ

EXERCISE

Table 3. Adjusted mean^a (95% confidence interval) of BMI and serum concentration of metabolic biomarkers in American adults by categories of weekly frequency of fast-food or pizza meals, NHANES 2007–2010

BMI or serum biomarker	Weekly frequency of fast-food or pizza meals				p ^b
	0 Time	1 Time	2–3 Times	≥ 4 Times	
BMI^c, kg m⁻²					
All (N = 8169)	27.5 (27.1, 27.8)	27.9 (27.6, 28.2)	28.9 (28.4, 29.4)	28.8 (28.3, 29.2)	< 0.0001
Men (n = 4002)	27.9 (27.4, 28.3)	28.0 (27.6, 28.4)	28.5 (28.0, 29.0)	28.6 (28.2, 29.0)	0.05
Women (n = 4167)	27.2 (26.8, 27.6)	27.7 (27.3, 28.1)	29.3 (28.6, 29.9)	29.0 (28.1, 29.8)	< 0.0001
Total cholesterol, mg dl ⁻¹ (N = 8236)	199 (197, 202)	198 (196, 200)	199 (196, 201)	198 (196, 201)	0.5
HDL-cholesterol^f, mg dl⁻¹					
All (n = 8236)	54 (53, 55)	53 (52, 54)	52 (51, 53)	51 (50, 52)	< 0.0001
Men (n = 4042)	48 (47, 49)	48 (47, 49)	48 (46, 49)	46 (45, 47)	0.003
Women (n = 4194)	60 (59, 61)	58 (57, 60)	56 (55, 57)	56 (54, 58)	0.001
LDL-cholesterol^d, mg dl⁻¹					
All (n = 3604)	113 (111, 116)	117 (113, 120)	113 (110, 116)	114 (110, 118)	0.6
< 50 Years (n = 2151)	107 (105, 110)	112 (109, 116)	111 (107, 114)	108 (104, 112)	0.8
≥ 50 Years (n = 1453)	123 (118, 129)	126 (121, 131)	118 (113, 123)	129 (122, 137)	0.5
Triglycerides, mg dl ⁻¹ (n = 3659)	103 (98, 109)	103 (99, 108)	110 (106, 115)	110 (104, 117)	0.2
Fasting glucose^e, mg dl⁻¹					
All (n = 3668)	99 (98, 100)	99 (98, 100)	99 (98, 100)	99 (98, 100)	0.5
Men (n = 1750)	102 (101, 104)	102 (101, 104)	101 (99, 102)	101 (99, 102)	0.1
Women (n = 1918)	97 (95, 98)	95 (94, 97)	97 (96, 99)	98 (96, 101)	0.2
Glycohemoglobin, % (N = 8234)	5.42 (5.39, 5.44)	5.39 (5.36, 5.42)	5.39 (5.36, 5.42)	5.40 (5.37, 5.44)	0.2

Abbreviations: BMI, body mass index; HDL, high-density lipoprotein; LDL, low-density lipoprotein; NHANES, National Health and Nutrition Examination Surveys. ^aAdjusted means were computed from multiple linear regression models with each biomarker as a continuous dependent variable. All biomarkers (except BMI, total- and HDL-cholesterol) were log-transformed for analysis; therefore, the back-transformed values for LDL-cholesterol, triglycerides, fasting glucose and glycohemoglobin are geometric means and their 95% confidence intervals. Independent variables included: frequency of fast-food meals (0, 1, 2–3 and ≥ 4 times), age (20–39, 40–59 and ≥ 60), sex, race/ethnicity (non-Hispanic white, non-Hispanic black, Mexican-American and other), poverty income ratio (≤ 1.3, > 1.3–3.5, ≥ 3.5 and unknown), years of education (< 12, 12, some college and ≥ college), serum cotinine (continuous), hours of fasting before phlebotomy, (continuous), physical activity (none, tertiles of MET minutes/week), alcohol-drinking status (never drinker, former drinker, current drinker and unknown). *N* refers to observations used in the regression model for each biomarker. ^b*P*-value for the Satterthwaite-adjusted *F* test for frequency of fast-food meals as a continuous variable. ^cSignificant interaction of fast-food meals with sex (*P*_{interaction} < 0.05; thus, the results are stratified by sex) ^dSignificant interaction of frequency of fast-food meals with age (*P*_{interaction} < 0.05); thus, the results are stratified by age categories.

INTRODUCTION

THE DATA SCIENCE WORKFLOW

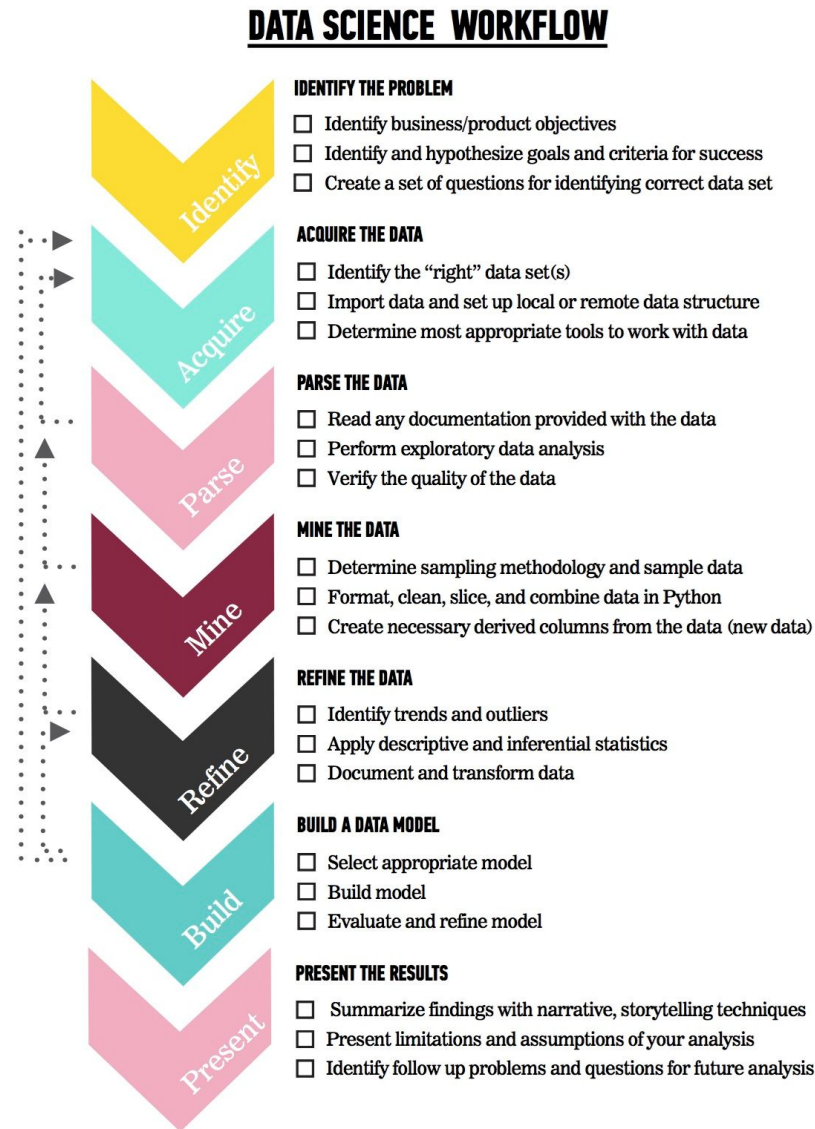
OVERVIEW OF THE DATA SCIENCE WORKFLOW

- A methodology for doing Data Science
- Similar to the scientific method
- Helps produce *reliable* and *reproducible* results
 - *Reliable*: Accurate findings
 - *Reproducible*: Others can follow your steps and get the same results

OVERVIEW OF THE DATA SCIENCE WORKFLOW

The steps:

1. Identify the problem
2. Acquire the data
3. Parse the data
4. Mine the data
5. Refine the data
6. Build a data model
7. Present the results



OVERVIEW OF THE 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

OVERVIEW OF THE DATA SCIENCE WORKFLOW



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

OVERVIEW OF THE DATA SCIENCE WORKFLOW



PARSE THE DATA

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

OVERVIEW OF THE DATA SCIENCE WORKFLOW



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)

OVERVIEW OF THE DATA SCIENCE WORKFLOW



REFINE THE DATA

- ☐ Identify trends and outliers
- ☐ Apply descriptive and inferential statistics
- ☐ Document and transform data

OVERVIEW OF THE DATA SCIENCE WORKFLOW



BUILD A DATA MODEL

- ☐ Select appropriate model
- ☐ Build model
- ☐ Evaluate and refine model

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100

OVERVIEW OF THE DATA SCIENCE WORKFLOW



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

EXAMPLE #1: PREDICTING PAYMENTS FRAUD

- **Problem:** Some advertisers attempt to use hacked accounts or stolen credit cards to run bad ads on Facebook
- **Goal:** Find the bad actors before their ads are delivered to users



EXAMPLE #2: LAUNCHING A NEW PRODUCT FEATURE

- › **Problem:** Users want more ways than just clicking “like” to respond to content in their news feed
- › **Goal:** Build and test additional “reactions”



GUIDED PRACTICE

DATA SCIENCE WORK FLOW

ACTIVITY: DATA SCIENCE WORKFLOW



EXERCISE

DIRECTIONS (25 minutes)

1. Divide into pairs, each located at a whiteboard.
2. **IDENTIFY:** Each group should develop 1 research question they would like to know about their classmates. Create a hypothesis to your question. Don't share your question yet! (5 minutes)
3. **ACQUIRE:** Rotate from group to group to collect data for your hypothesis. Have other students write or tally their answers on the whiteboard. (10 minutes)
4. **PRESENT:** Communicate the results of your analysis to the class. (10 minutes)
 - a. Create a narrative to summarize your findings.
 - b. Provide a basic visualization for easy comprehension.
 - c. Choose one student to present for the group.

DELIVERABLE

Presentation of the results

WELCOME TO DATA SCIENCE

EXIT TICKET

DON'T FORGET TO FILL OUT YOUR EXIT TICKET

CONCLUSION

REVIEW

CONCLUSION

- 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?