

WELCOME TO DATA SCIENCE

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

PRE-WORK

PRE-WORK REVIEW

- ▶ Define basic data types used in object-oriented programming
- ▶ Recall the Python syntax for lists, dictionaries, and functions

DATA SCIENCE

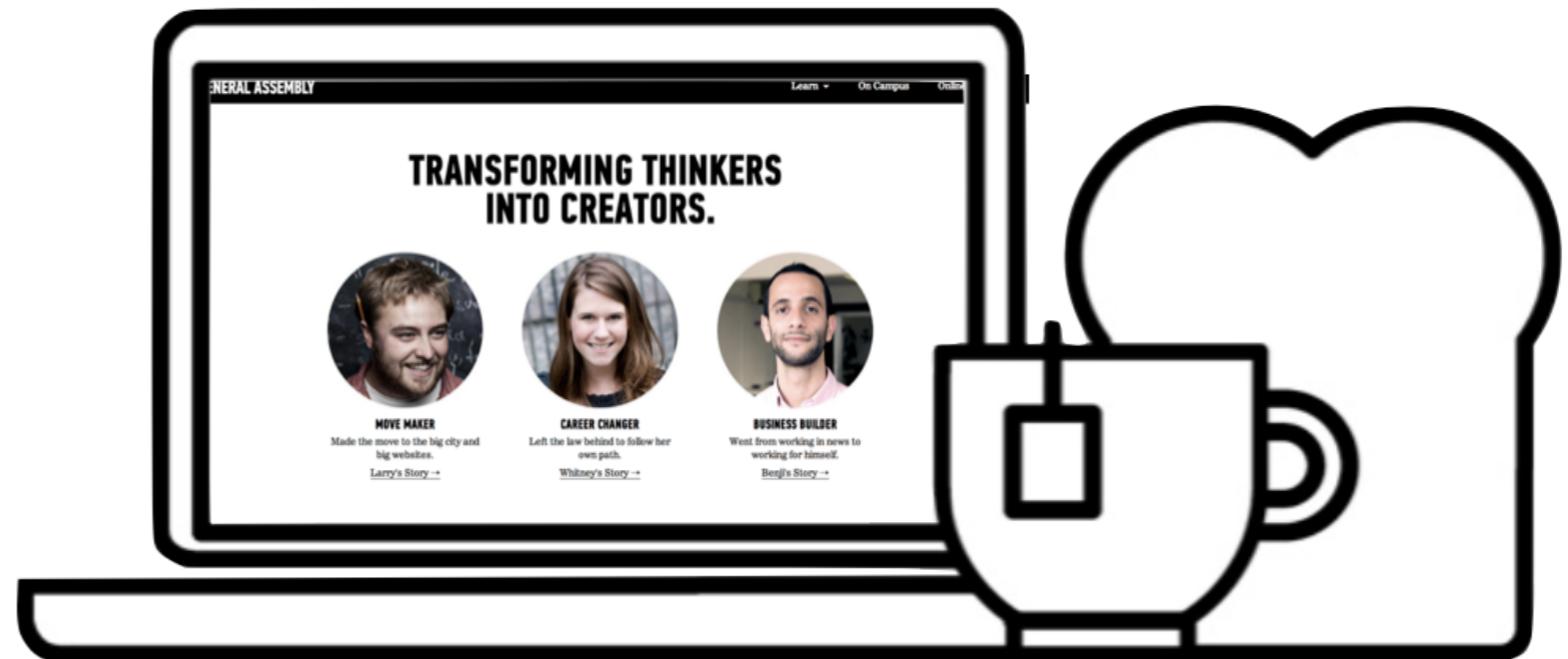
WELCOME TO GA!

WELCOME TO GA!

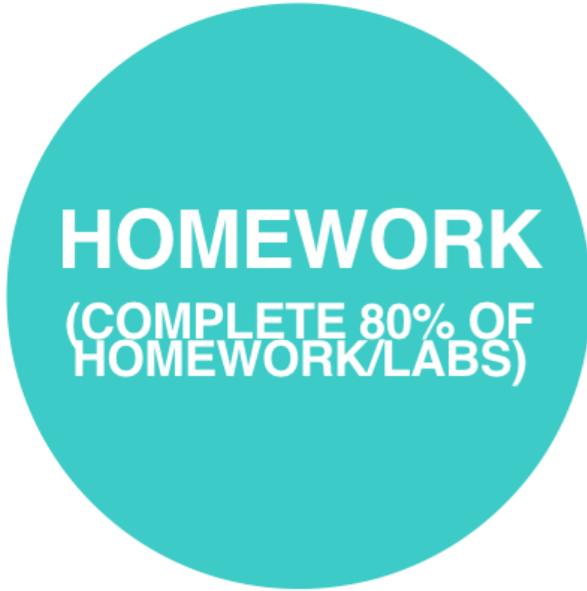
- General Assembly is a global community of individuals empowered to pursue the work we love.
- General Assembly's mission is to build our community by transforming millions of thinkers into creators.

FEEDBACK/SUPPORT

- ▶ Access to EIRs: office hours, in class support
- ▶ Exit Tickets
- ▶ Mid-Course Feedback
- ▶ End of Course Feedback



GA GRADUATION REQUIREMENTS



HOMEWORK
(COMPLETE 80% OF
HOMEWORK/LABS)



ATTENDANCE
(MISS NO MORE THAN 2
CLASSES)



**FINAL
PROJECT**



**COMMUNITY
ENGAGEMENT**
PARTICIPATION +
FEEDBACK

FOREVER AND EVER



**BUILD
YOUR
NETWORK**

It's not just about altruism, your network is your most valuable asset



**FIND
OPPORTU
NITIES**

Alumni have started companies together and recruited other alumni to join their teams



**13,000+
STRONG**

You're part of the alumni community forever



PERKS!
15% OFF CLAASSES
AND WORKSHOPS, \$500 TUITION CREDIT

We can't wait to have you back on campus

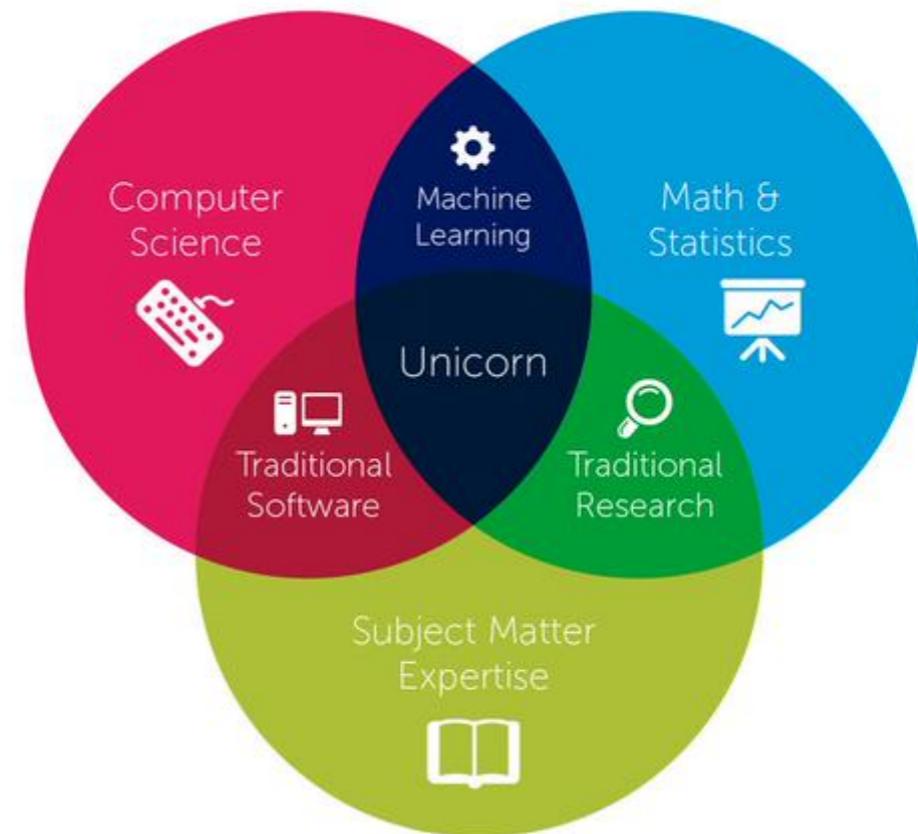
INTRODUCTION

WHAT IS DATA SCIENCE?

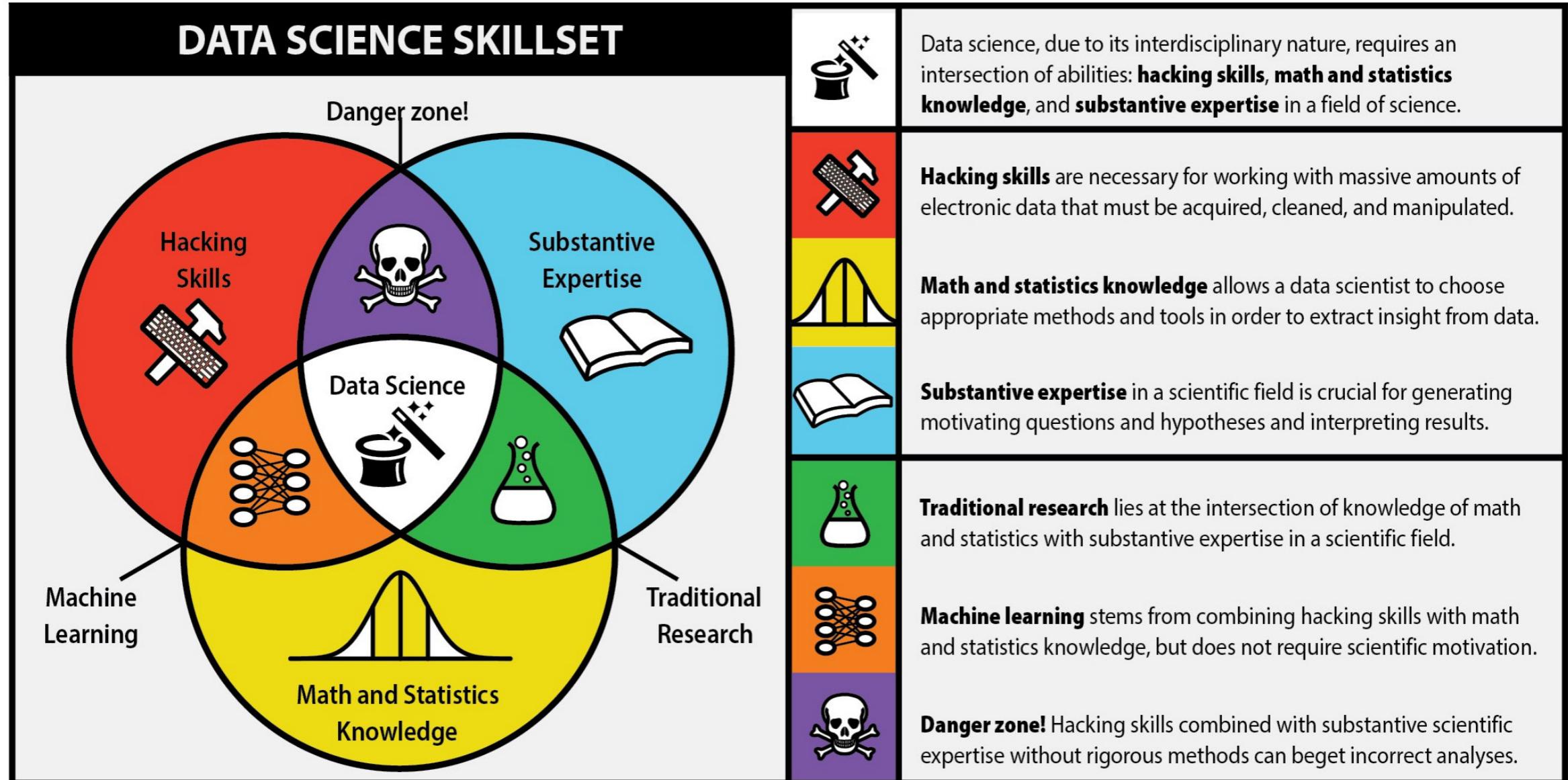
WHAT IS DATA SCIENCE?

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

Data Science



WHAT IS DATA SCIENCE?



WHO USES DATA SCIENCE?

NETFLIX

amazon.com®

Google



FiveThirtyEight



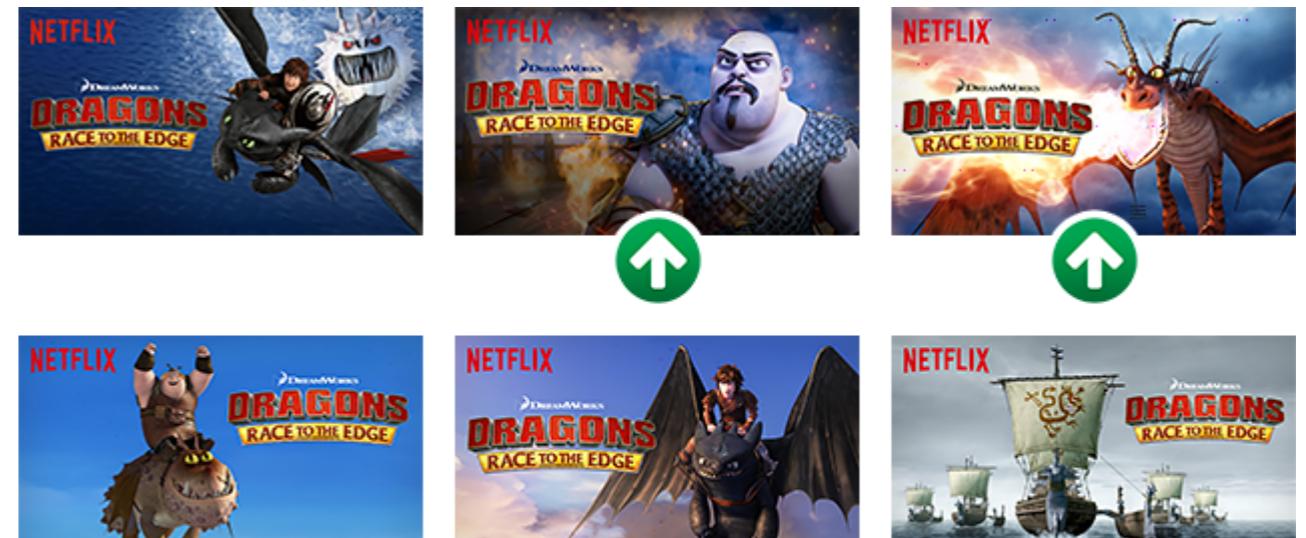
WHO USES DATA SCIENCE?

The Switch

Red light, green light: New Audis will predict the time until that stoplight turns green [Article](#)



Innovations



Netflix reveals what images hook viewers on new shows [Article](#)

WHO USES DATA SCIENCE?

► Can you think of others?

WHAT IS DATA SCIENCE?

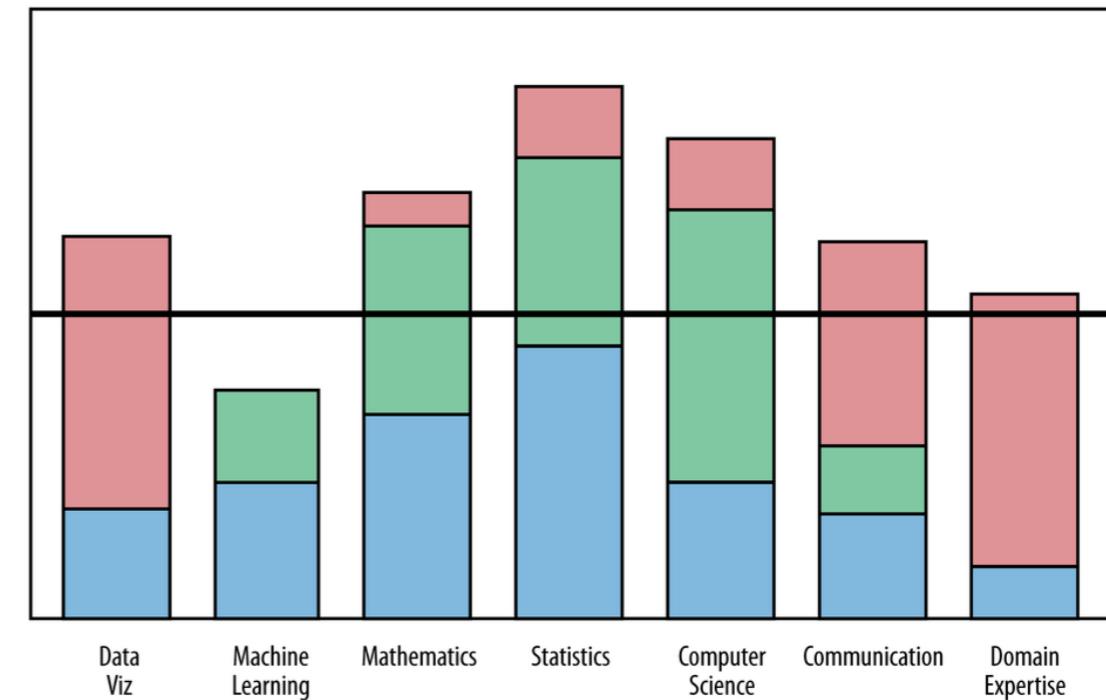
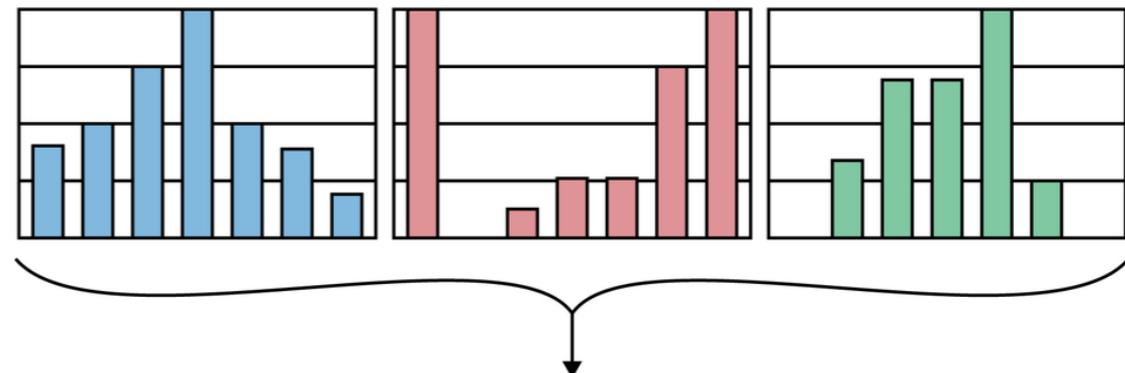
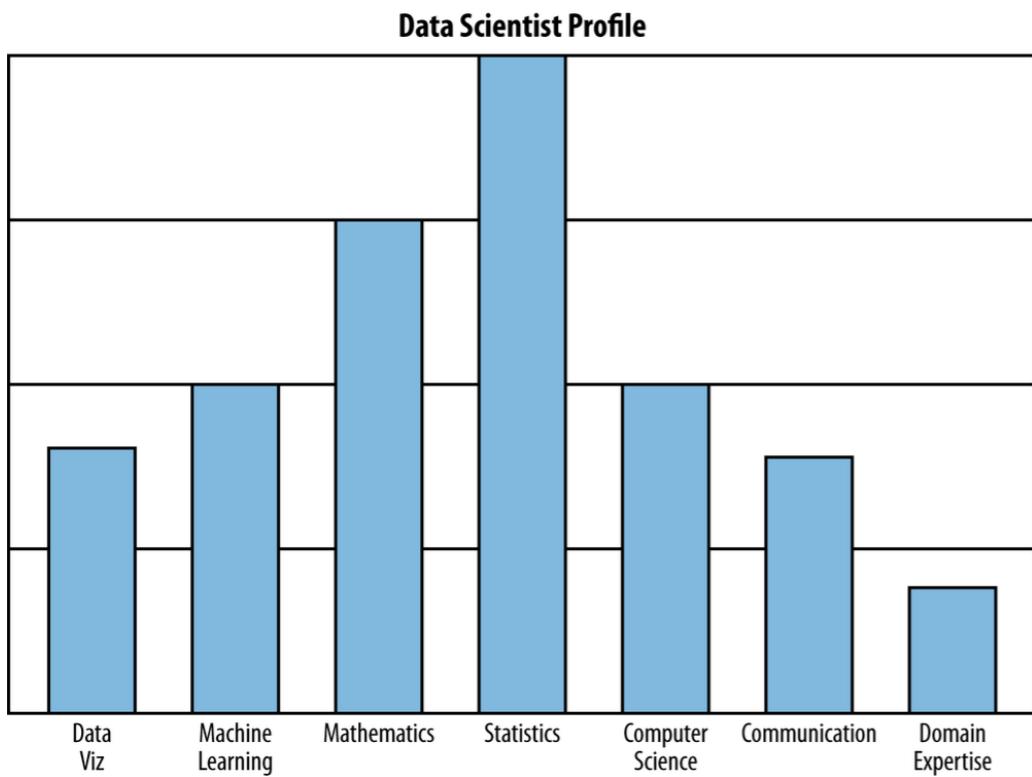


INTRODUCTION

Data Science Roles

WHAT IS DATA SCIENCE?

No one person can be the perfect data scientist, so **we need teams**.



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

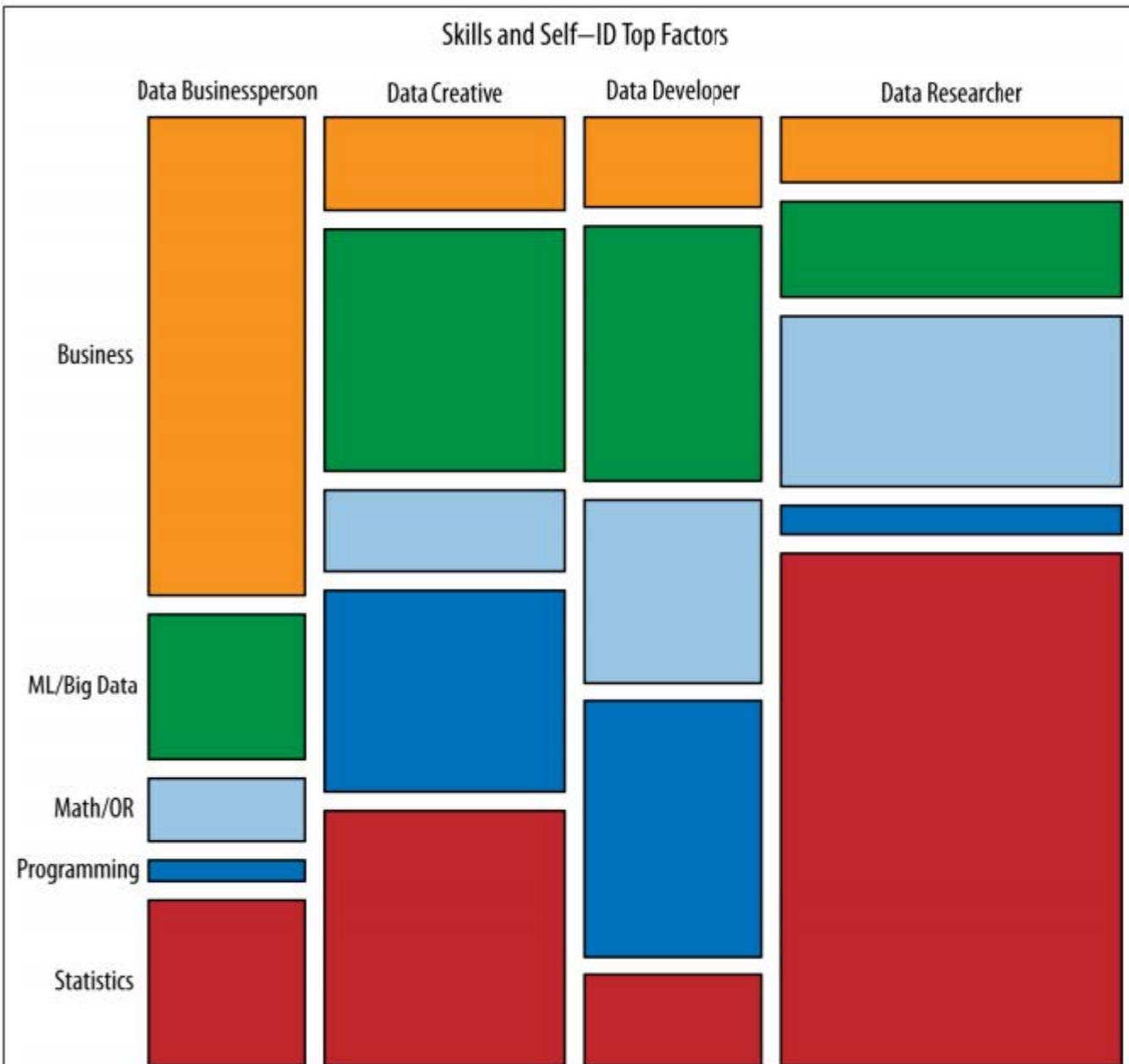
WHAT ARE THE ROLES IN DATA SCIENCE?

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

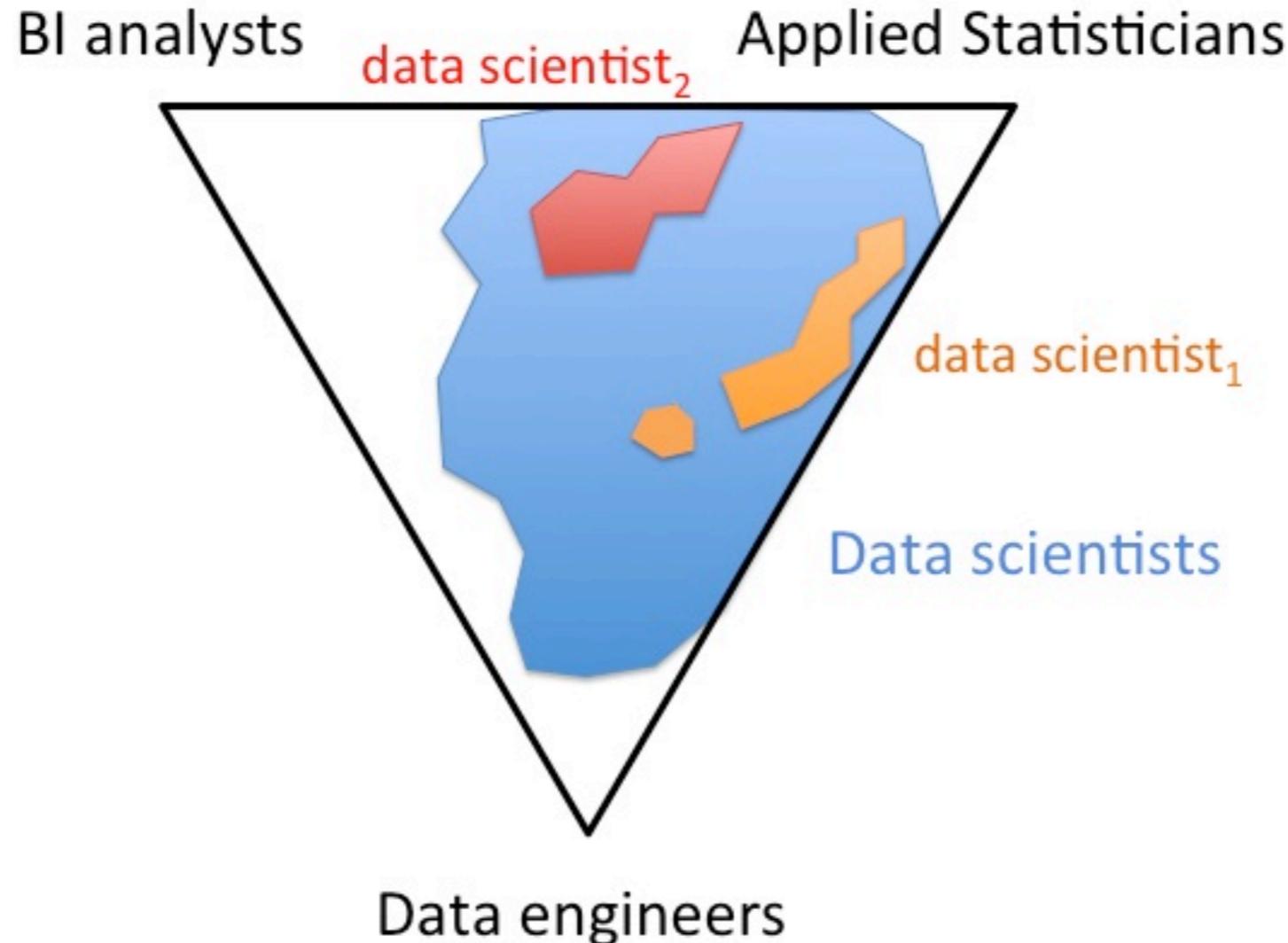
Business	ML / Big Data	Math / OR	Programming	Statistics
Product Development Business	Unstructured Data Structured Data Machine Learning Big and Distributed Data	Optimization Math Graphical Models Bayesian / Monte Carlo Statistics Algorithms Simulation	Systems Administration Back End Programming Front End Programming	Visualization Temporal Statistics Surveys and Marketing Spatial Statistics Science Data Manipulation Classical Statistics

WHAT ARE THE ROLES IN DATA SCIENCE?

- ▶ These roles prioritize different skill sets.
- ▶ However, all roles involve some part of each skillset.
- ▶ Where are your strengths and weaknesses?



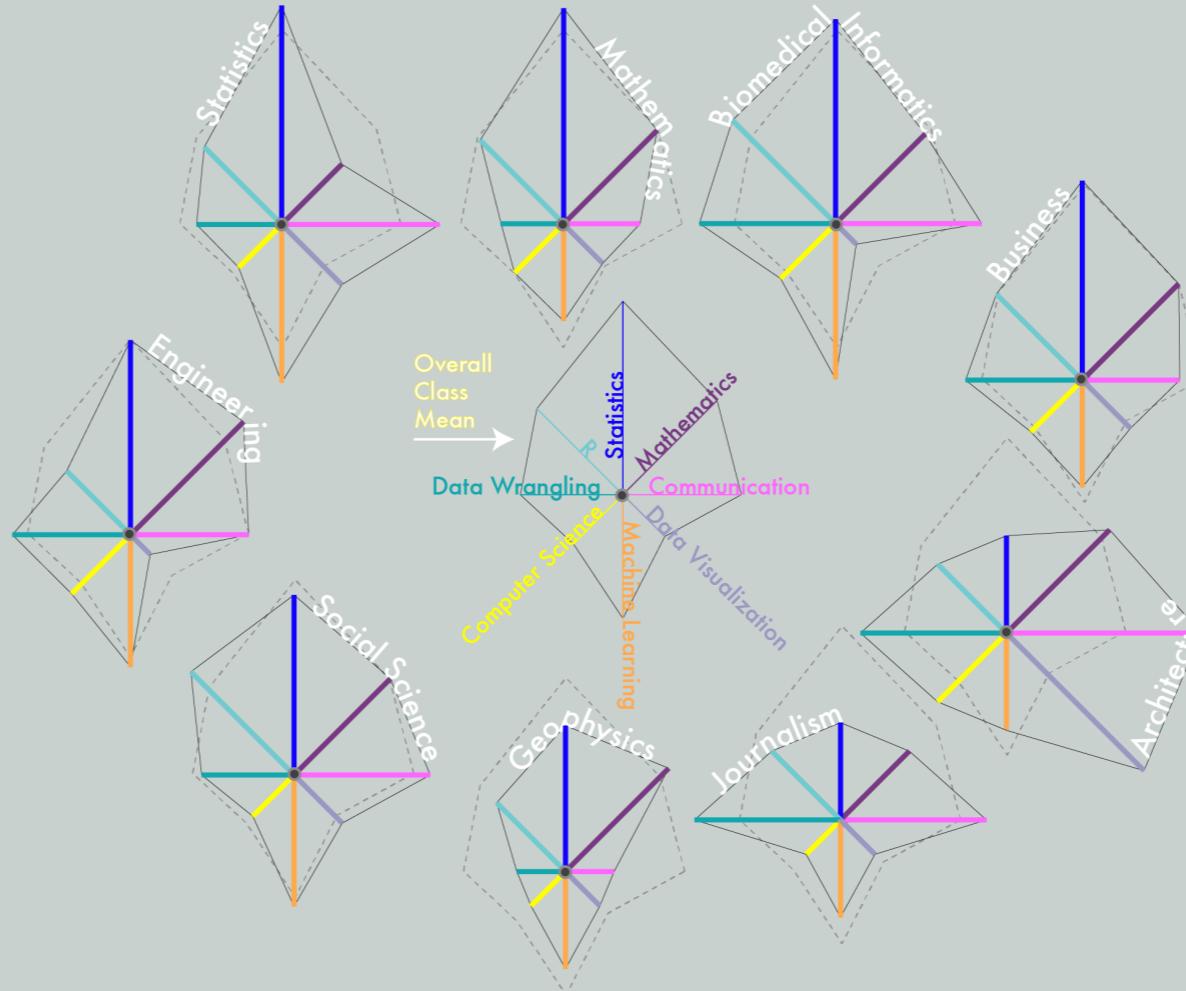
WHAT ARE THE ROLES IN DATA SCIENCE?



WHAT ARE THE ROLES IN DATA SCIENCE

The Stars of Data Science

Students in Columbia's Introduction to Data Science course came from across the academic spectrum. Their skills are presented here in star charts with spokes representing their skill levels* across the data science skillset: R, statistics, mathematics, communication, data visualization, machine learning, computer science, and data wrangling. In addition to hovering in the center, the star chart of the overall class mean underlies each academic domain, so you can see students from each academic domain relative to the rest of the class. How would you compose your own intergalactic data science team?



*Skills were assessed by a survey written and administered by a subset of students in the class.

INTRODUCTION

THE DATA SCIENCE WORKFLOW

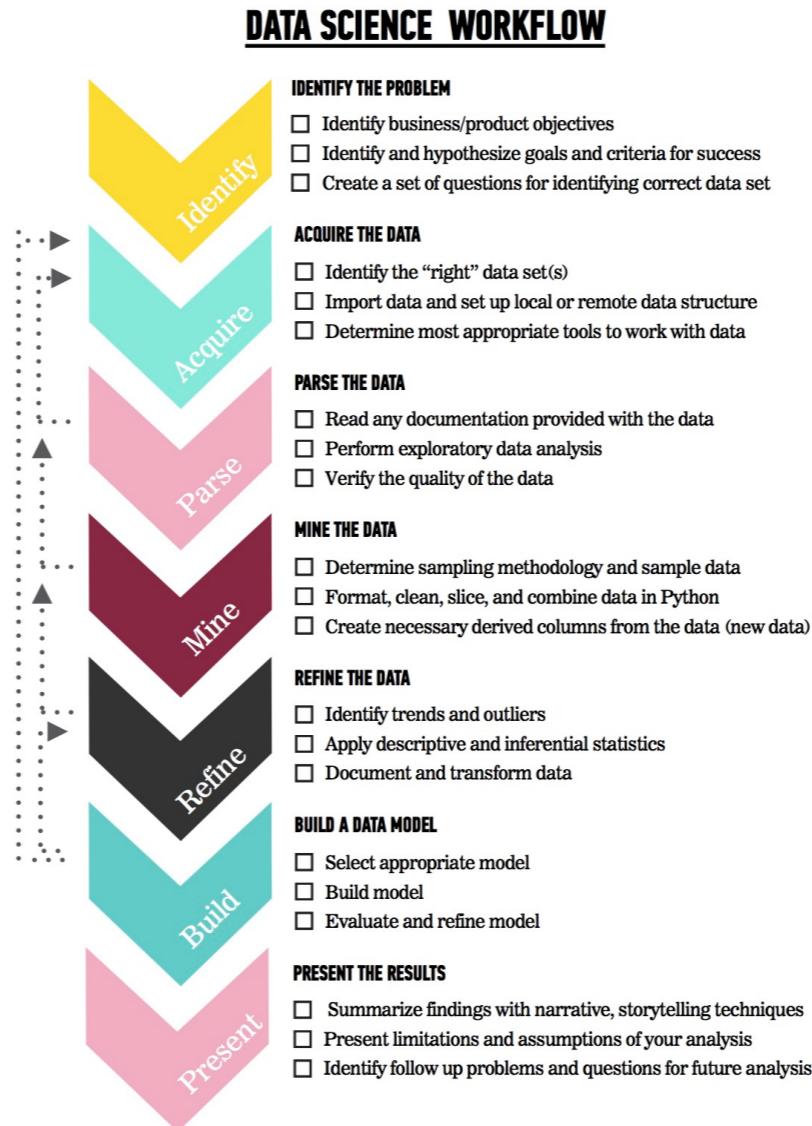
OVERVIEW OF THE DATA SCIENCE WORKFLOW

- ▶ 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



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

OVERVIEW OF THE DATA SCIENCE WORKFLOW

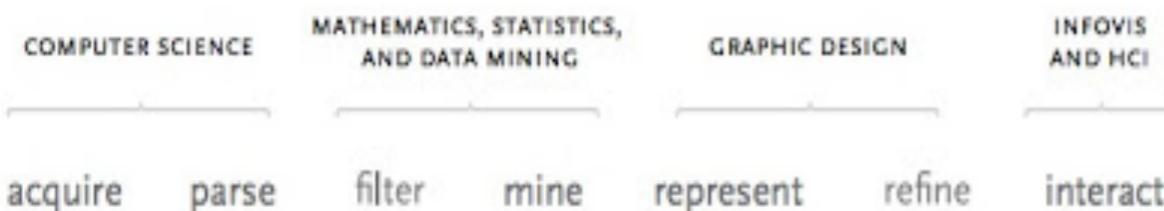


PRESENT THE RESULTS

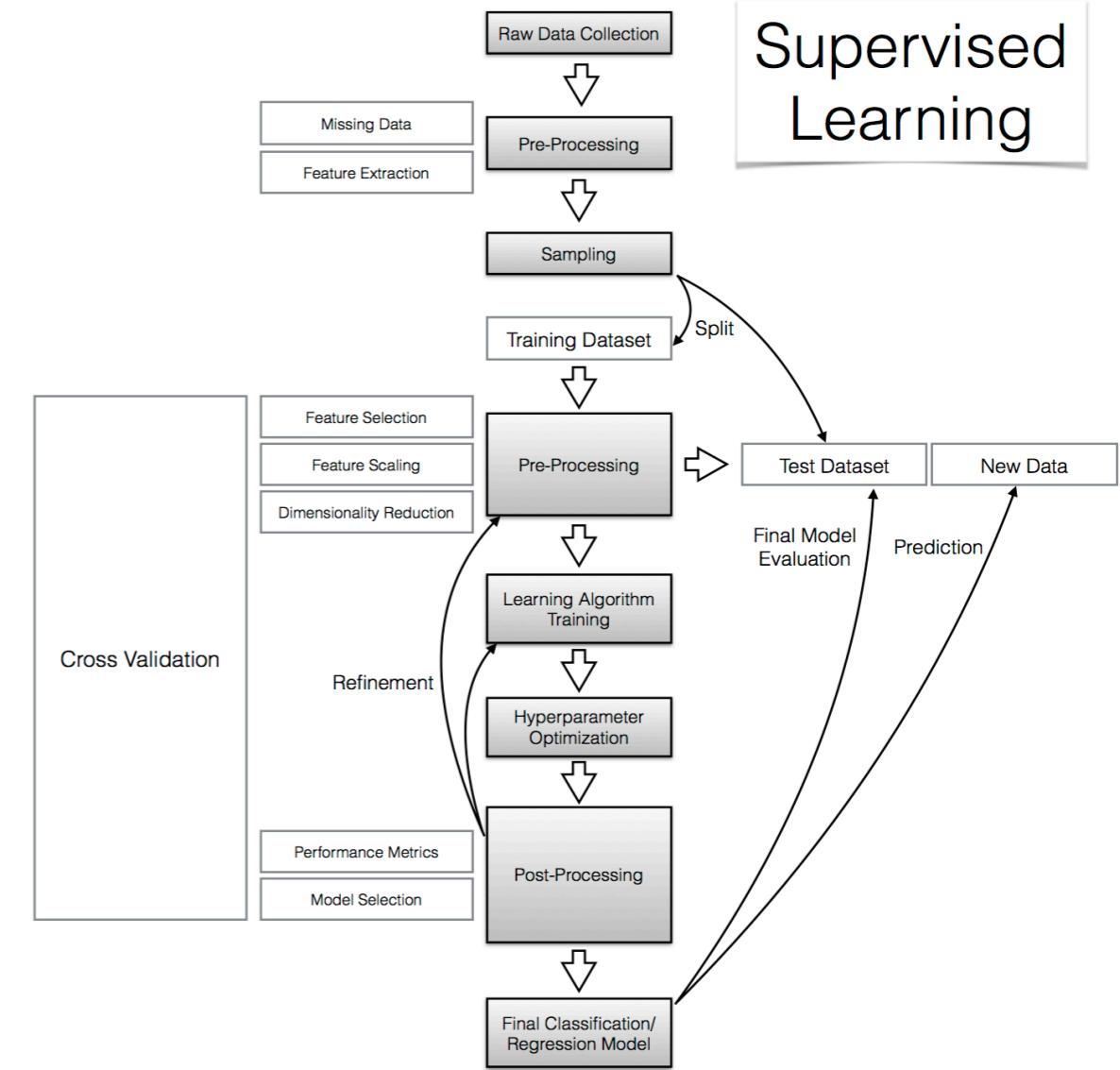
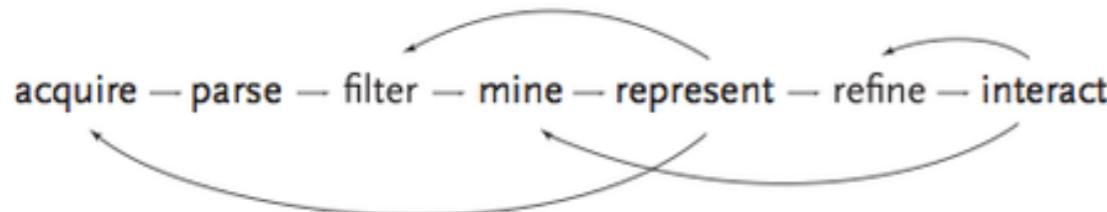
- Summarize findings with narrative, storytelling techniques
- Present limitations and assumptions of your analysis
- Identify follow up problems and questions for future analysis

OVERVIEW OF THE DATA SCIENCE WORKFLOW

Interdisciplinary



Recursion



Sebastian Raschka 2014

This work is licensed under a Creative Commons Attribution 4.0 International License.

GUIDED PRACTICE

DATA SCIENCE WORK FLOW EXERCISE

ACTIVITY: DATA SCIENCE WORKFLOW



DIRECTIONS (25 minutes)

1. Divide into 4 groups, 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

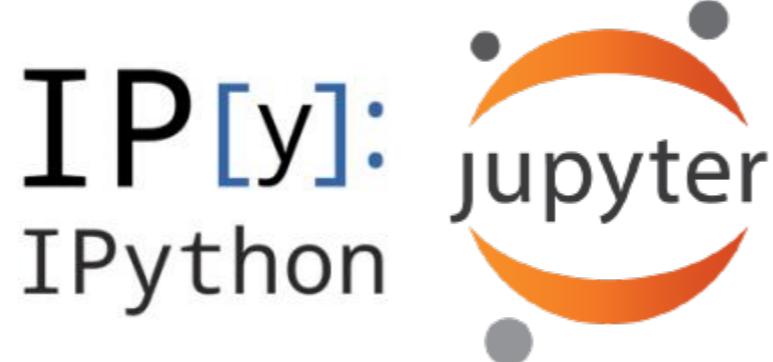
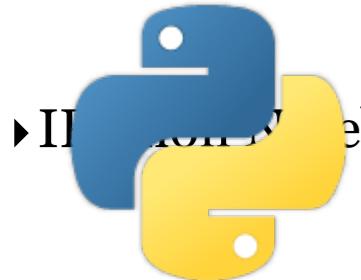
Presentation of the results

DEMO

ENVIRONMENT SETUP

DEV ENVIRONMENT SETUP

- ▶ A Brief Intro to Our Environment



DEV ENVIRONMENT SETUP

- ▶ Brief intro of tools
- ▶ Environment setup
 - ▶ Create a Github account
 - ▶ Install Python 2.7 and Anaconda
 - ▶ Practice Python syntax, Terminal commands, and Pandas
- ▶ iPython Notebook test and Python review

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?

DATA SCIENCE

**BEFORE NEXT
CLASS**

WELCOME TO DATA SCIENCE

Q & A

WELCOME TO DATA SCIENCE

EXIT TICKET

**DON'T FORGET TO FILL OUT YOUR EXIT
TICKET**