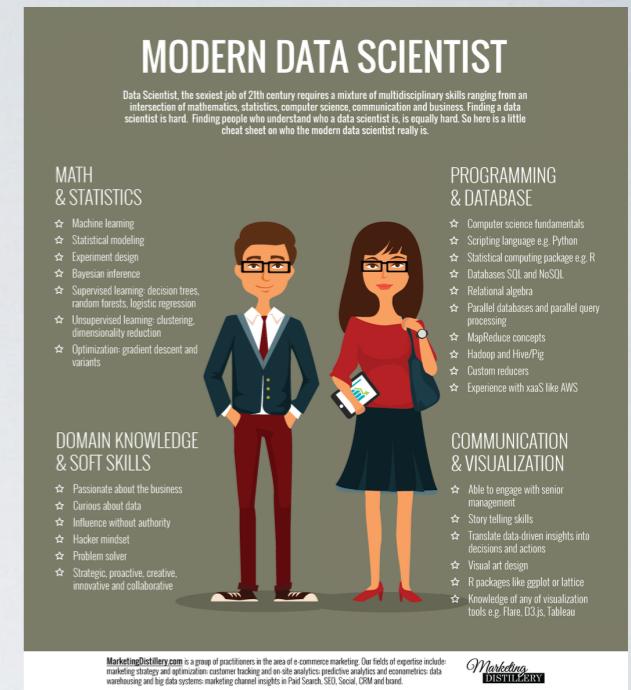
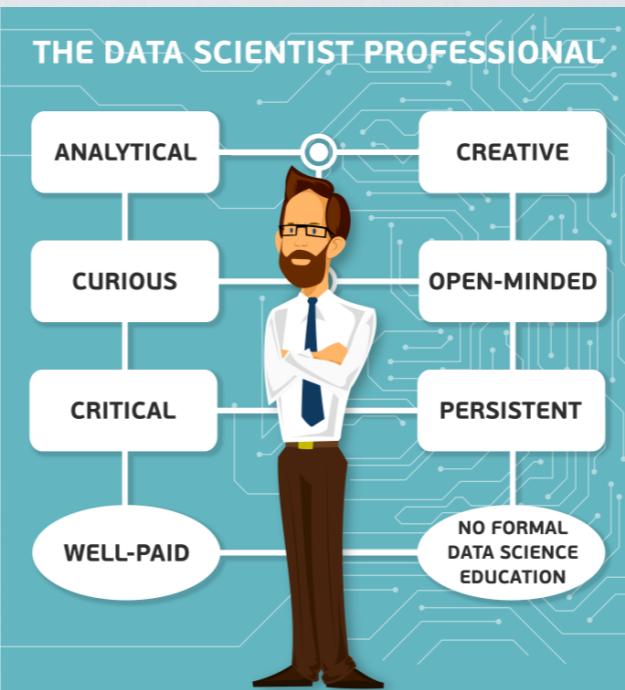
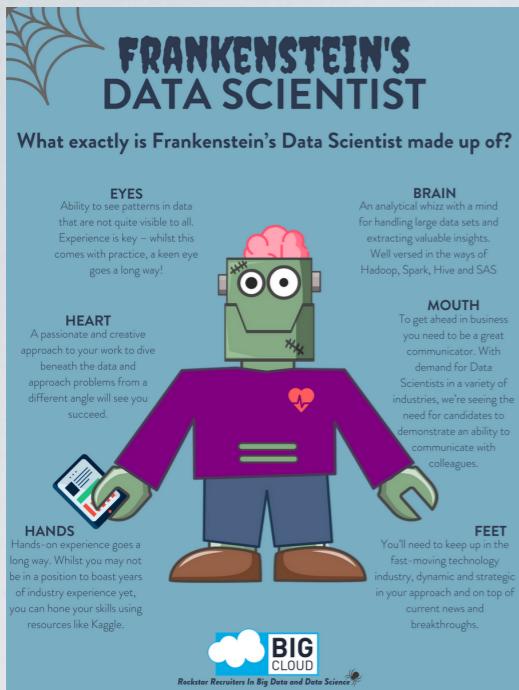




EFFECTIVE DATA SCIENCE

Jan 16, 2019 - Columbus Data Science

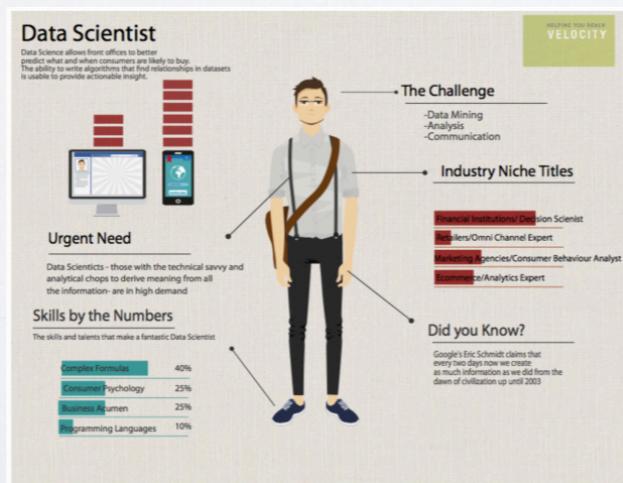
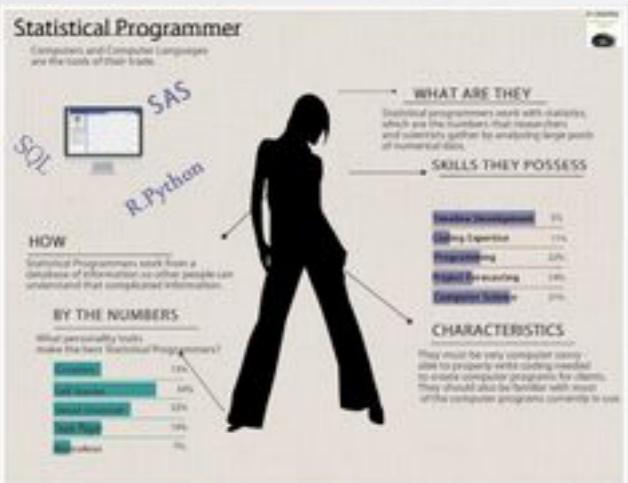
WHAT IS(N'T) DATA SCIENCE



Who are Data Scientists?

(As defined by Infographics on the first page of a google search)

6/9 - Male
 4/9 - Bearded
 8/9 - Human
 2/9 - Wear ties



1 Data Scientist



4.8 / 5
Job Score

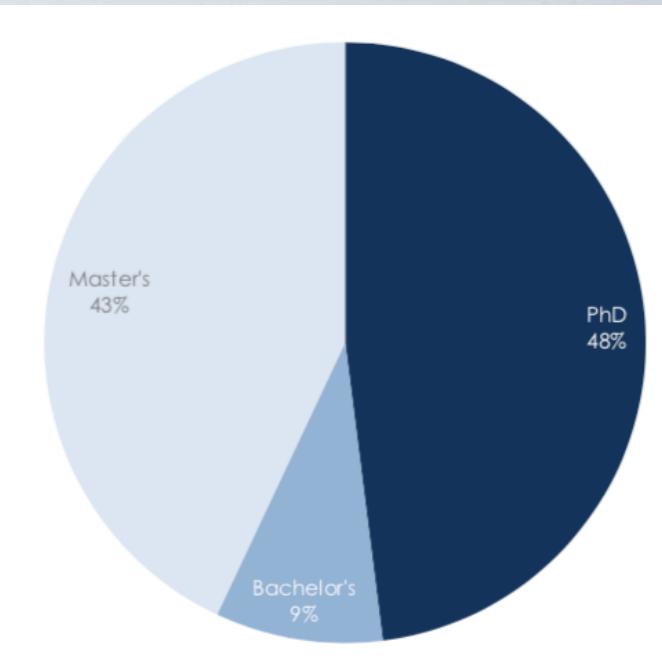
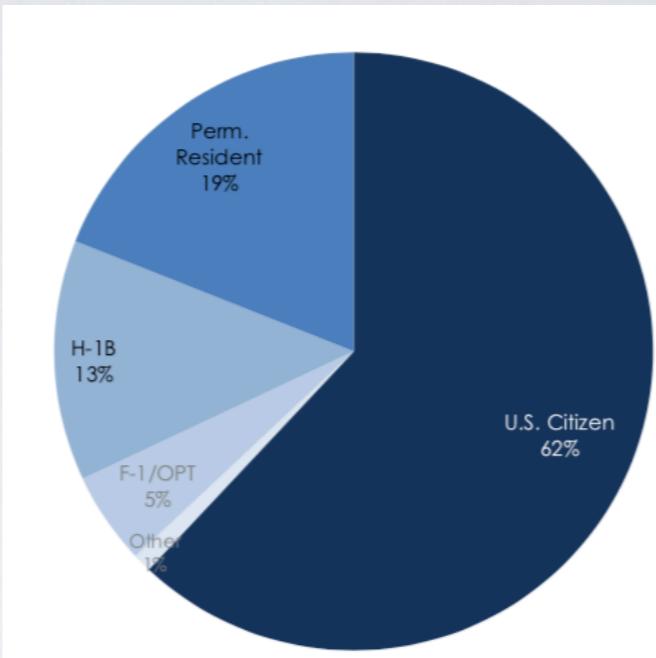
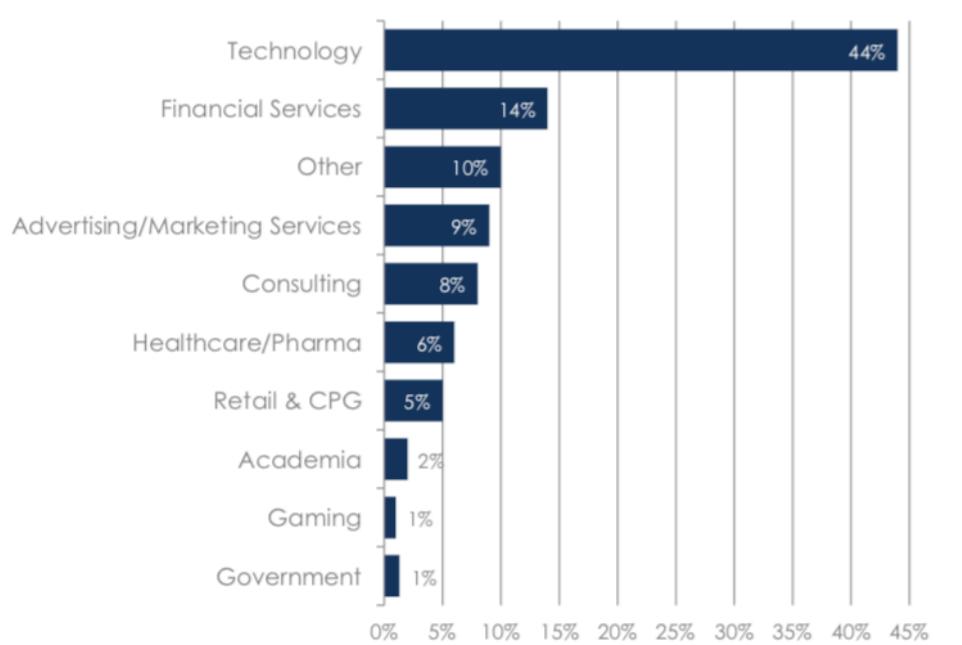
\$110,000
Median Base Salary

4.2 / 5
Job Satisfaction

4,524
Job Openings

[View Jobs](#)





Who are Data Scientists?

(As defined by a Burtch Works Survey)

91% - Have an advanced degree

40% on West Coast

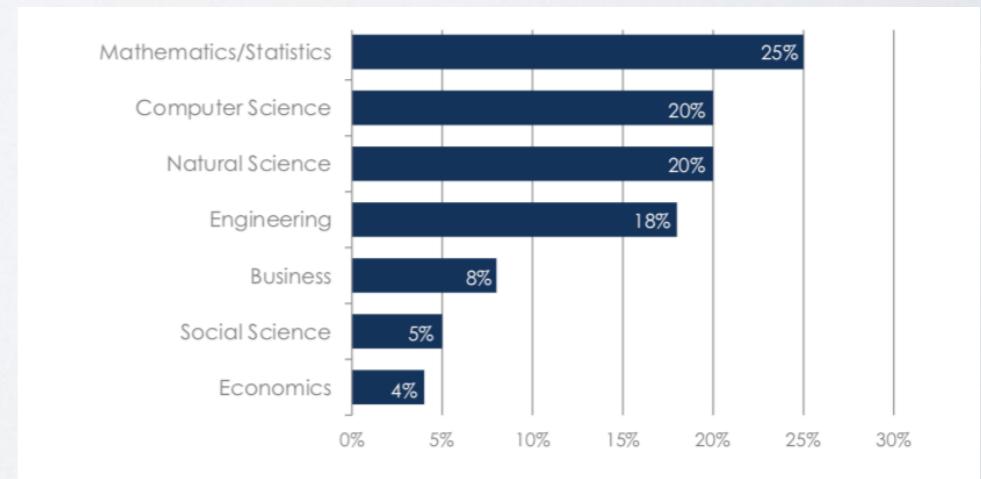
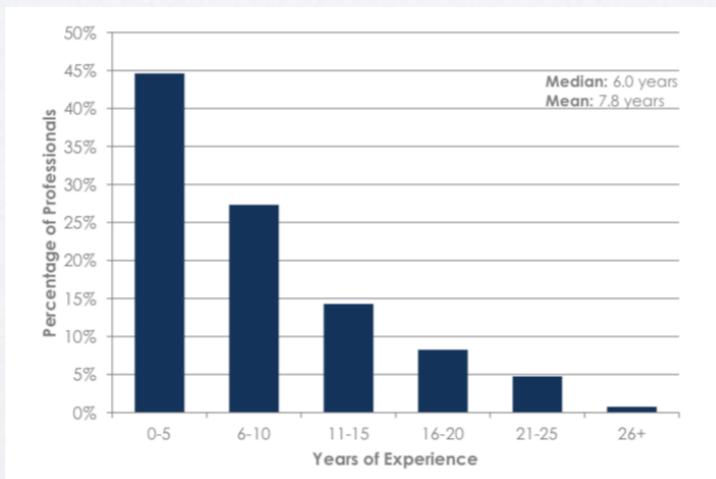
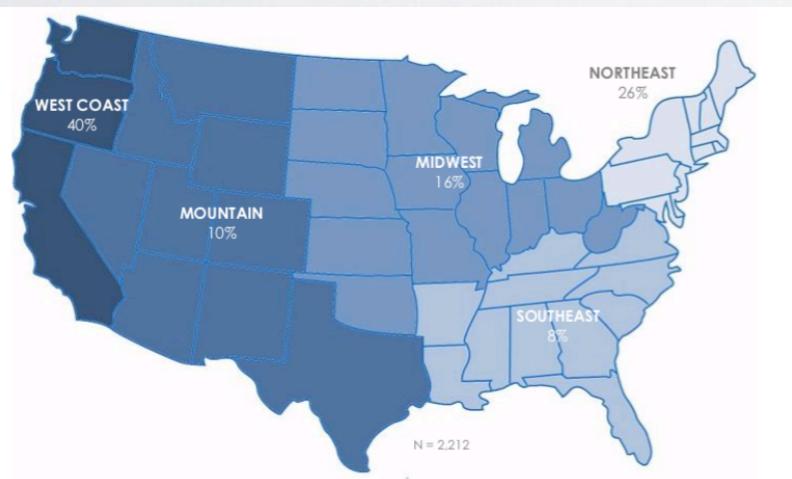
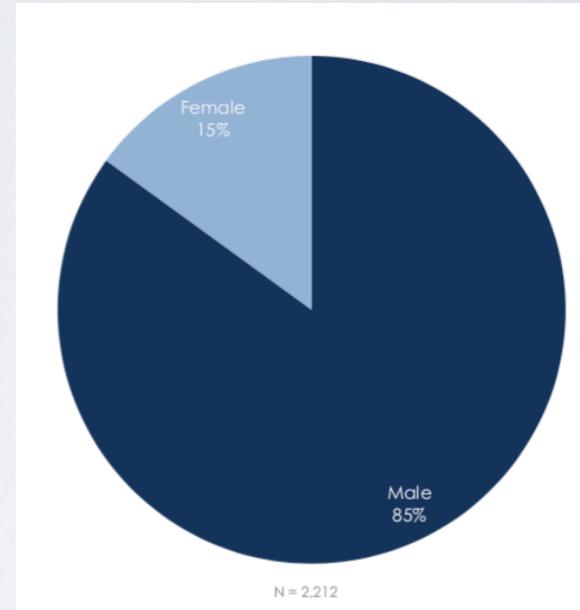
26% in Northeast

44% in Tech industry

85% are male

Median years of experience - 6

Note: 4/9 faces in Butch Works slide deck has a beard



WHAT IS DATA SCIENCE?

CCD DATA: THE *GOOD*, THE *BAD*, AND THE *UGLY*

PHILIP MASSEY and GEORGE H. JACOBY

Kitt Peak National Observatory, NOAO, P.O. Box 26732, Tucson, AZ
85726-6732

“Some days the magic works. Some days it doesn’t.”

Little Big Man’s grandfather.

INTRODUCTION

We will describe three kinds of data: *good data*, *bad data*, and *ugly data* (also known as “truth”). The first of these, *good data*, actually doesn’t exist, but rather is an ideal to strive for in terms of what you would like your CCD data to be like. *Bad data* consists of data that you are better off throwing out—something went wrong; you either didn’t notice or couldn’t do anything about it, but just throw it out. *Ugly data* is what you get from *real* CCDs, and is just fine—just a little bit, well, ugly. It’s our plan to help you distinguish between the last two cases.

WHAT IS DATA SCIENCE?

Think

Pair

Share

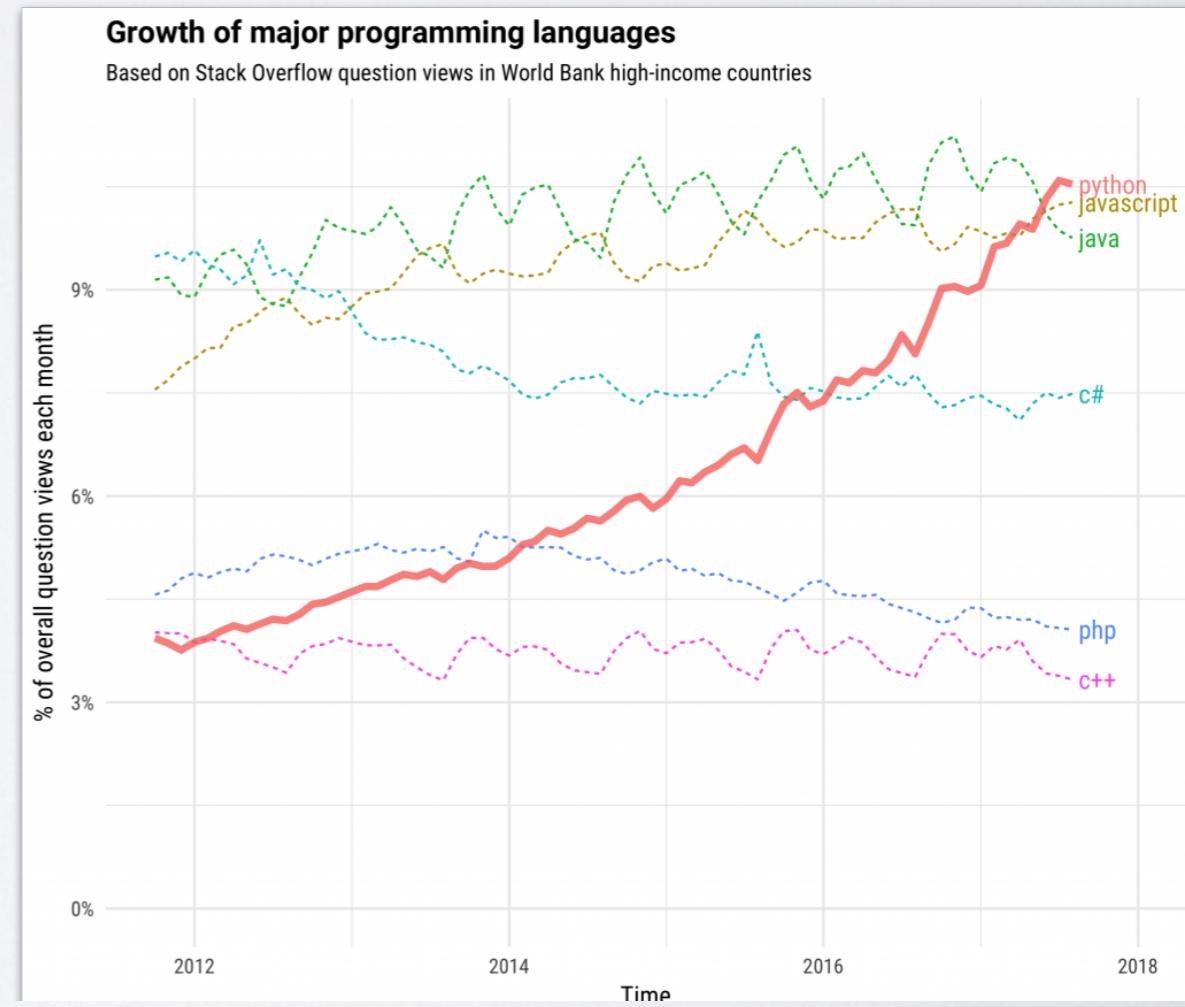
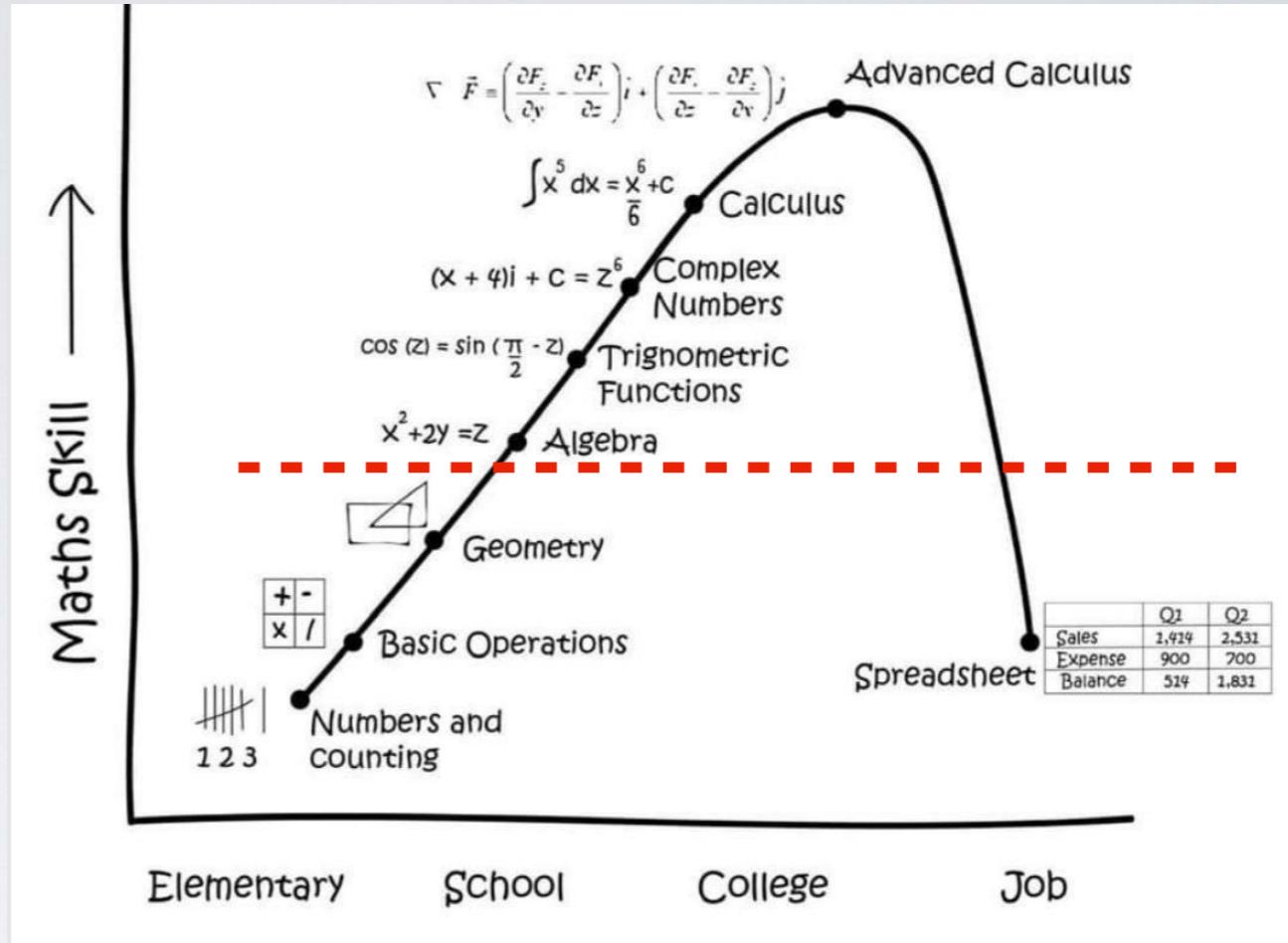
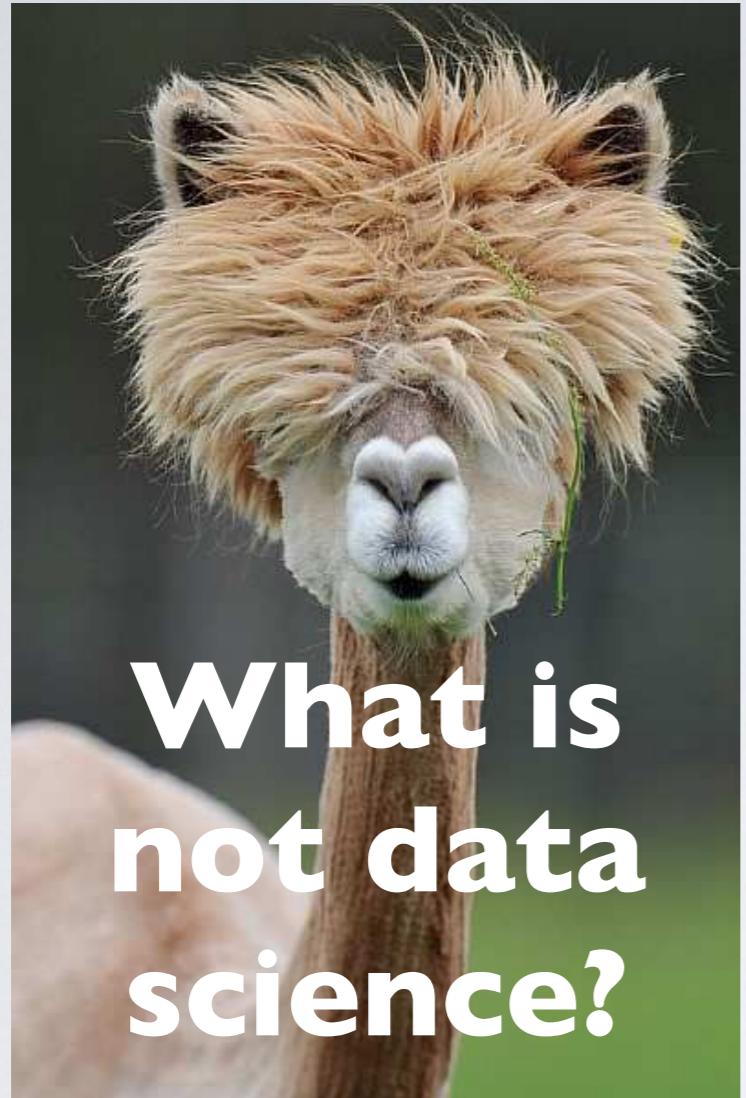


Source: IBM® BigInsights™

WHAT IS DATA SCIENCE?

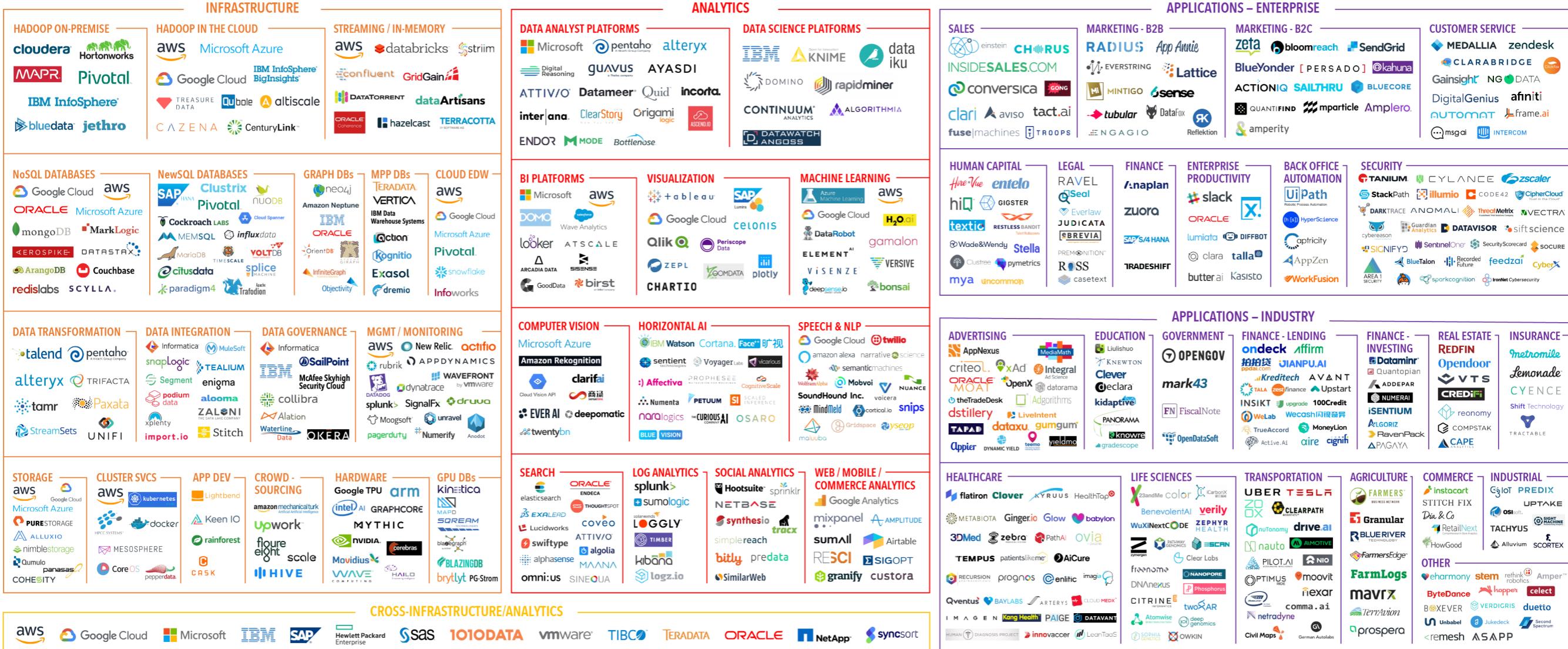
- Critical Thinking
- Problem Solving
- Computer Programming
- Analytics
- Machine Learning
- Stuff and stuff





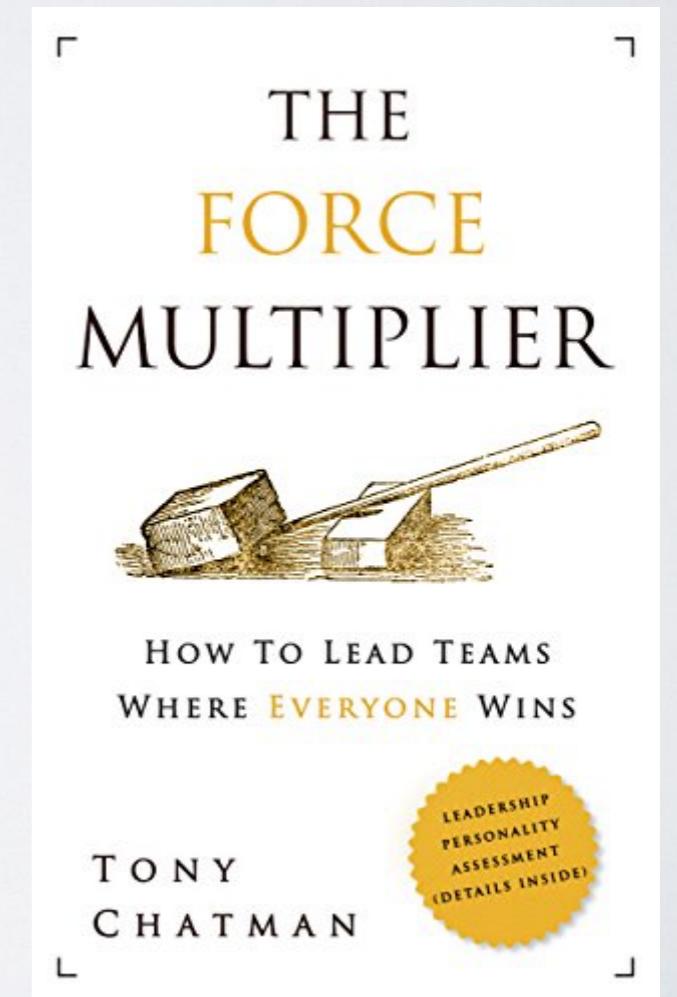
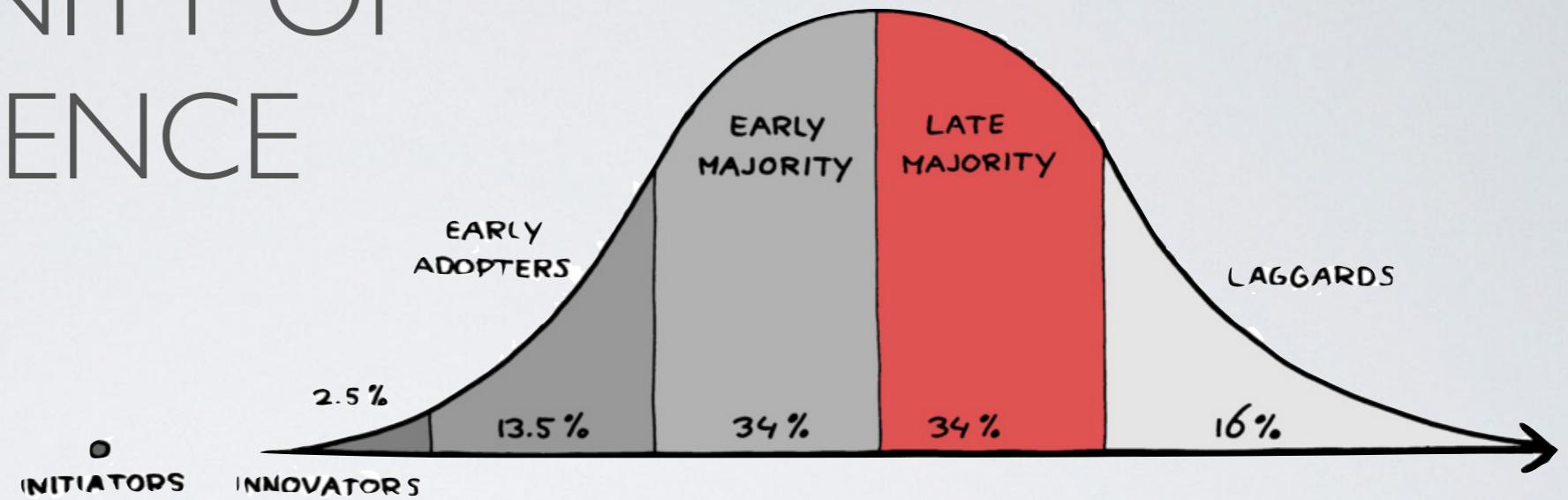
HOW DOES DS BENEFIT
FROM COMMUNITY

BIG DATA & AI LANDSCAPE 2018



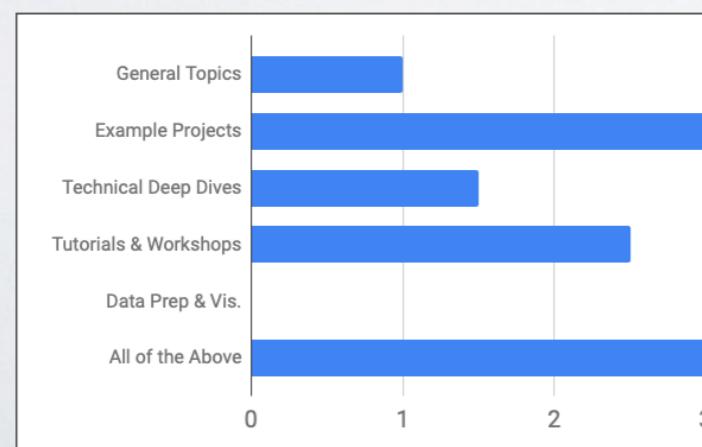
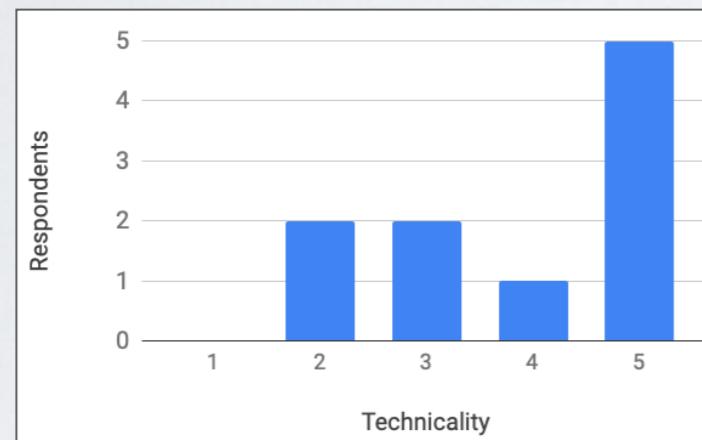
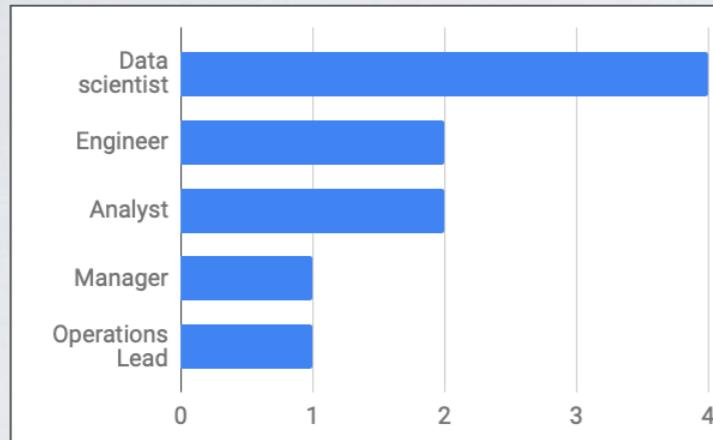
A COMMUNITY OF DATA SCIENCE

- Network
- Learn stuff
- Showcase your skills
- Get answers to questions
- Discuss new technologies
- Explain where you are stuck
- Brainstorm
- Force Multiplier



WHAT DOES OUR
COMMUNITY WANT/NEED

COLUMBUS DATA SCIENCE POLL



C-Bus Data Science Audience Survey

What's your current role?

- Data scientist
- Engineer
- Analyst
- Manager
- Other: _____

How technical are you?

- 1 2 3 4 5
- I'm interested and want to learn more
- I code on the reg

What sort of talks are interesting to you?

- What is data science? / Building data science teams
- Data science projects / a day in the life
- Technical deep dives
- Tutorials & workshops
- Data exploration & visualization
- Other: _____

Anything else we should know?

Your answer

A COMMUNITY OF DATA SCIENCE

- Network
- Share projects
- Share tools
- Showcase your skills
- Discuss new technologies
- Practice presenting
- Discuss new technologies
- Argue over the use of tabs or spaces

MODERN DATA SCIENTIST

Data Scientist, the sexiest job of 21th 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



DOMAIN KNOWLEDGE & SOFT SKILLS

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

PROGRAMMING & DATABASE

- ★ Computer science fundamentals
- ★ Scripting language e.g. Python
- ★ Statistical computing package 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

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

“This is our circus lets train those monkeys to dance!”

– Socrates