# Data Science in the Wild

## Lecture 1: Introduction

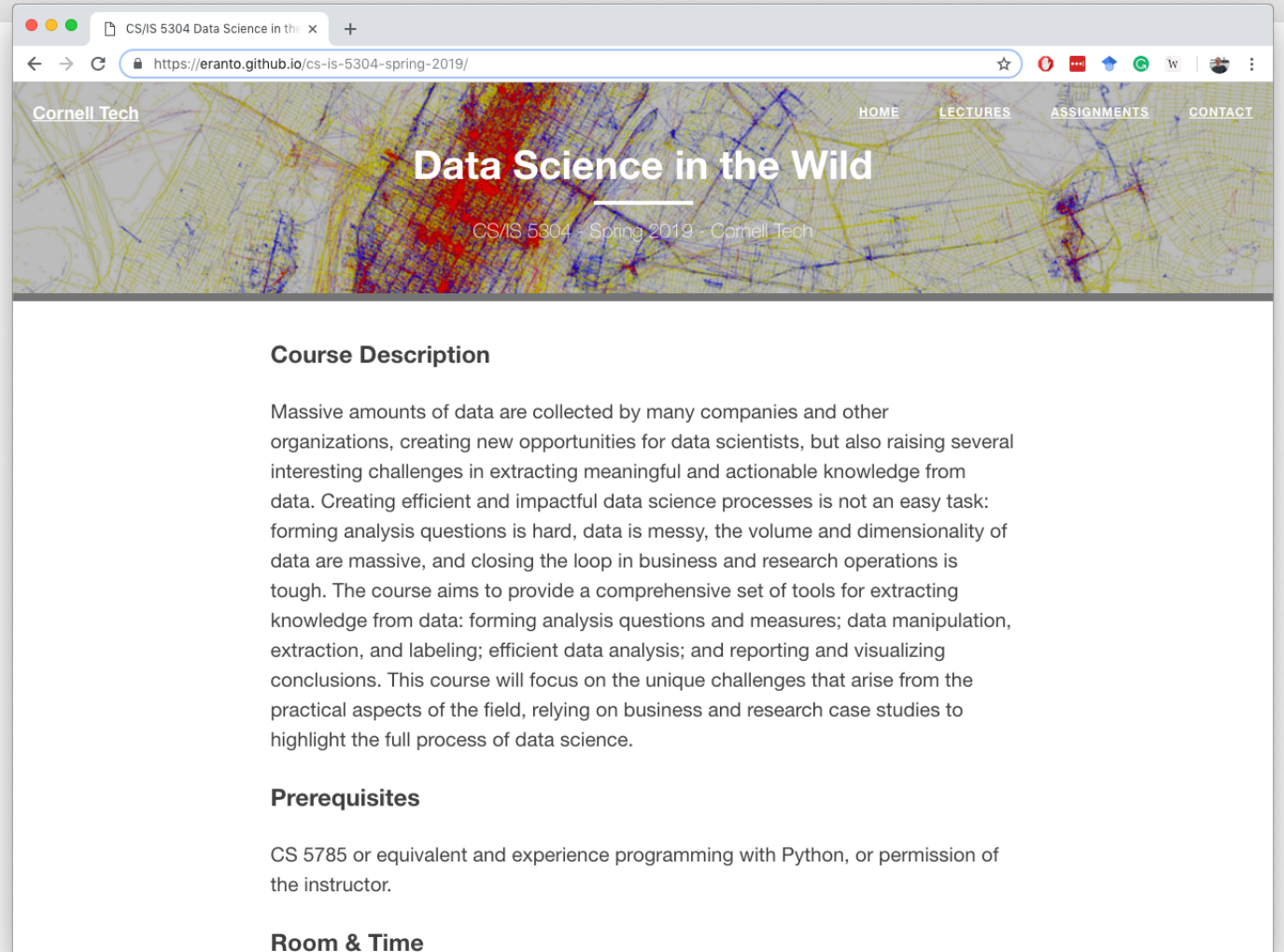Eran Toch

CORNELL TECH

# Agenda

1. About the Course

2. The Data Explosion

3. Data Science Capabilities

4. The scientific method

# Resources

- Website: https://eranto.github.io/cs5304-spring2019/
- Slack: wild-data-science.slack.com

# Prof. Eran Toch

Visiting associate Professor at Cornell Tech
Faculty, Tel Aviv University
etoch@cornell.edu
Twitter: @erant
http://toch.tau.ac.il

# Mr. David Rimshnick

- Cornell OR alum, BS 2005, MEng 2006
  - Research on logistics problems (airline crew scheduling, vehicle routing)
- Spent career in data science and analytics in healthcare industry
  - ZS Associates
  - Novo Nordisk (biopharma company)
  - Pfizer (biopharma company)
- Currently Principal at Boston Consulting Group
  - Part of BCG Gamma, sub-organization devoted to advanced AI and ML applications

david.rimshnick@gmail.com

# Team

- TA:
  - Zekun Hao

- Graders:
  - Summer Shi
  - Seye Bankole
  - Svava Kristinsdottir
  - Mohit Chawla

# Timetable

Please let us know
about absence days
due to religious
holidays

| Lecture | Date | Lecture | Assignments |
|---|---|---|---|
| 1 | Jan 23, 2019 | Introduction to Data Science | |
| 2 | Jan 28, 2019 | Extract, Transform and Load | |
| 3 | Jan 30, 2019 | Cleaning and Labeling Data | Assignment 1 Due |
| 4 | Feb 4, 2019 | Learning from Unbalanced Data | |
| 5 | Feb 6, 2019 | Data labeling and Data Labelers | |
| 6 | Feb 11, 2019 | Analyzing Experiments | Assignment 2 Due |
| 7 | Feb 13, 2019 | Statistical Analysis of Experiments | |
| 8 | Feb 18, 2019 | Bias and Quality Measures | |
| 9 | Feb 20, 2019 | Data-Based Simulation / Impact Analysis | |
| 10 | Feb 25, 2019 | FEBRUARY BREAK | |
| 11 | Feb 27, 2019 | Big Data Tools for Data Science | |
| 12 | Mar 4, 2019 | Learning in Distributed Processing | Assignment 3 Due |
| 13 | Mar 6, 2019 | Programming Cache-Based Distributed Processing | |
| 14 | Mar 11, 2019 | Technical Topic - Hands on With Spark/PySpark | |
| 15 | Mar 13, 2019 | Company Presentation - Deep Learning for Drug Discovery (Stephen Ra, Pfizer) | Assignment 4 Due |
| 16 | Mar 18, 2019 | Preliminary exam | |
| 17 | Mar 20, 2019 | Deep Sequence Learning | |
| 18 | Mar 25, 2019 | Data Visualization | |
| 19 | Mar 27, 2019 | Deep Recommendation Systems | Project Part 1 Due |
| 20 | Apr 1, 2019 | SPRING BREAK | |
| 21 | Apr 3, 2019 | SPRING BREAK | |
| 22 | Apr 8 | Background: Reinforcement Learning | |
| 23 | Apr 10 | Reinforcement Learning | |
| 24 | Apr 15, 2019 | Guest Lecture (Samar Deen?) | |
| 25 | Apr 17, 2019 | Causality versus Correlation / Causal Effects | Project Part 2 Due |
| 26 | Apr 22, 2019 | LIME and Model Explainability | |
| 27 | Apr 24, 2019 | Communicating Results | |
| 28 | Apr 29, 2019 | Ethics of Data Science | |
| 29 | May 1, 2019 | Final Projects in Class | Final Project Due |
| 30 | May 6, 2019 | Final Projects in Class | Final Project Due |

# Grade Breakdown

- Home assignments (30%)
- Final project (30%)
- Preliminary exam (20%) - in class
- Final exam (20%) - take home

# Assignments

- 4 home assignments
  - Each with programming and a written exercise
  - Each students has a total of one slip day
- The officially supported programming language is Python
  - But you are welcome to work on your assignments using other languages
- You can use well-known libraries but cite them.
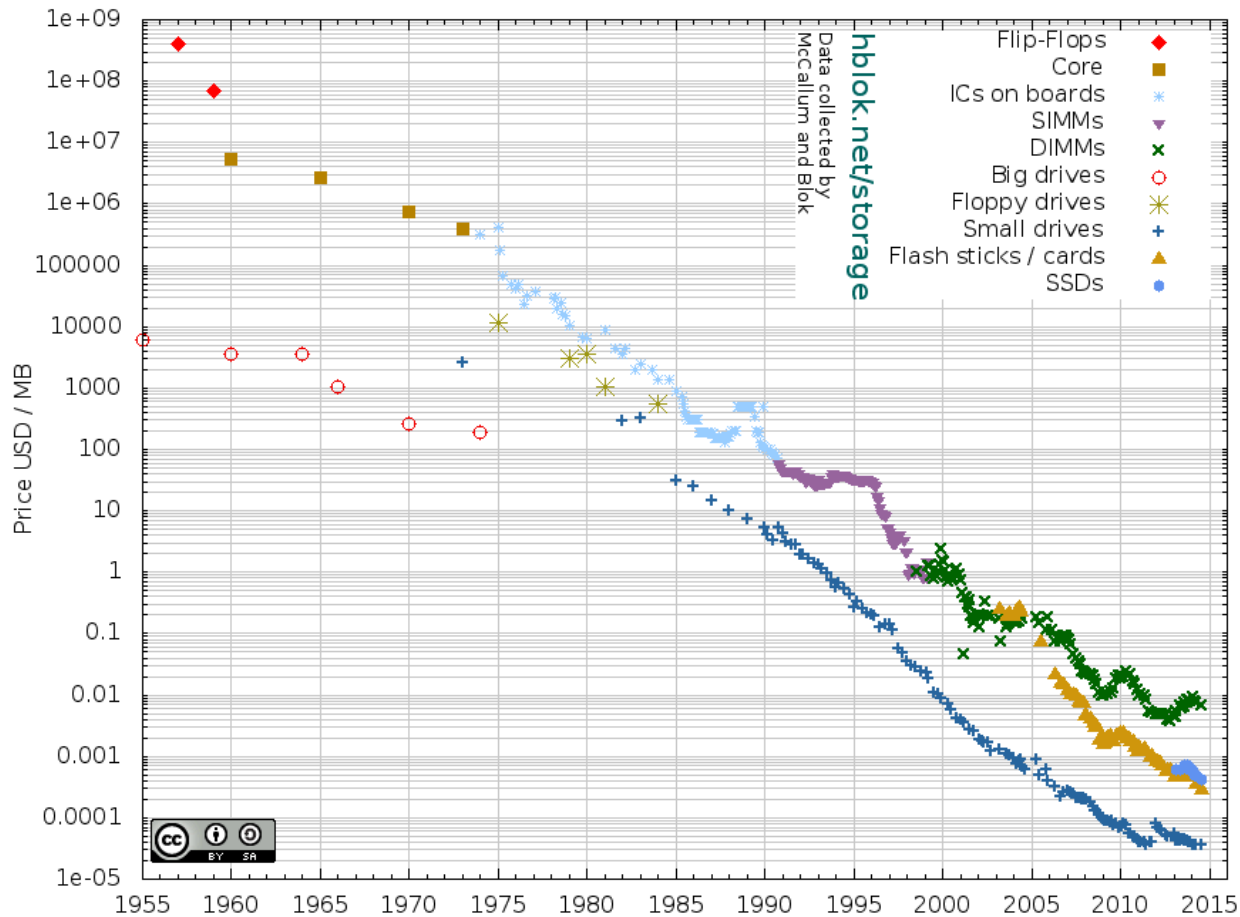- You are encouraged to work in groups of 2 students.

# Bibliography

**The books are not required for the course, but they can be of interest to students.**

1. Foster Provost and Tom Fawcett, Data Science for Business: What You Need to Know about Data Mining and Data-Analytic Thinking, O'Reilly Media; 1st edition (2013)

2. Jake VanderPlas, Python Data Science Handbook, O'Reilly Media; 1 edition (2016) - Free book

3. Russell Jurney, Agile Data Science 2.0: Building Full-Stack Data Analytics Applications with Spark, O'Reilly Media; 1st edition (2017).

4. A. Rajaraman, J. Leskovec and J. Ullman, Mining of Massive Datasets, Cambridge University Press, 3rd version

# Data Storage Prices



Historical Cost of Computer Memory and Storage

3.75 Megabyte

1 Terrabyte

# How do we make decisions?



According to HiPPO

(highest paid person's opinion)
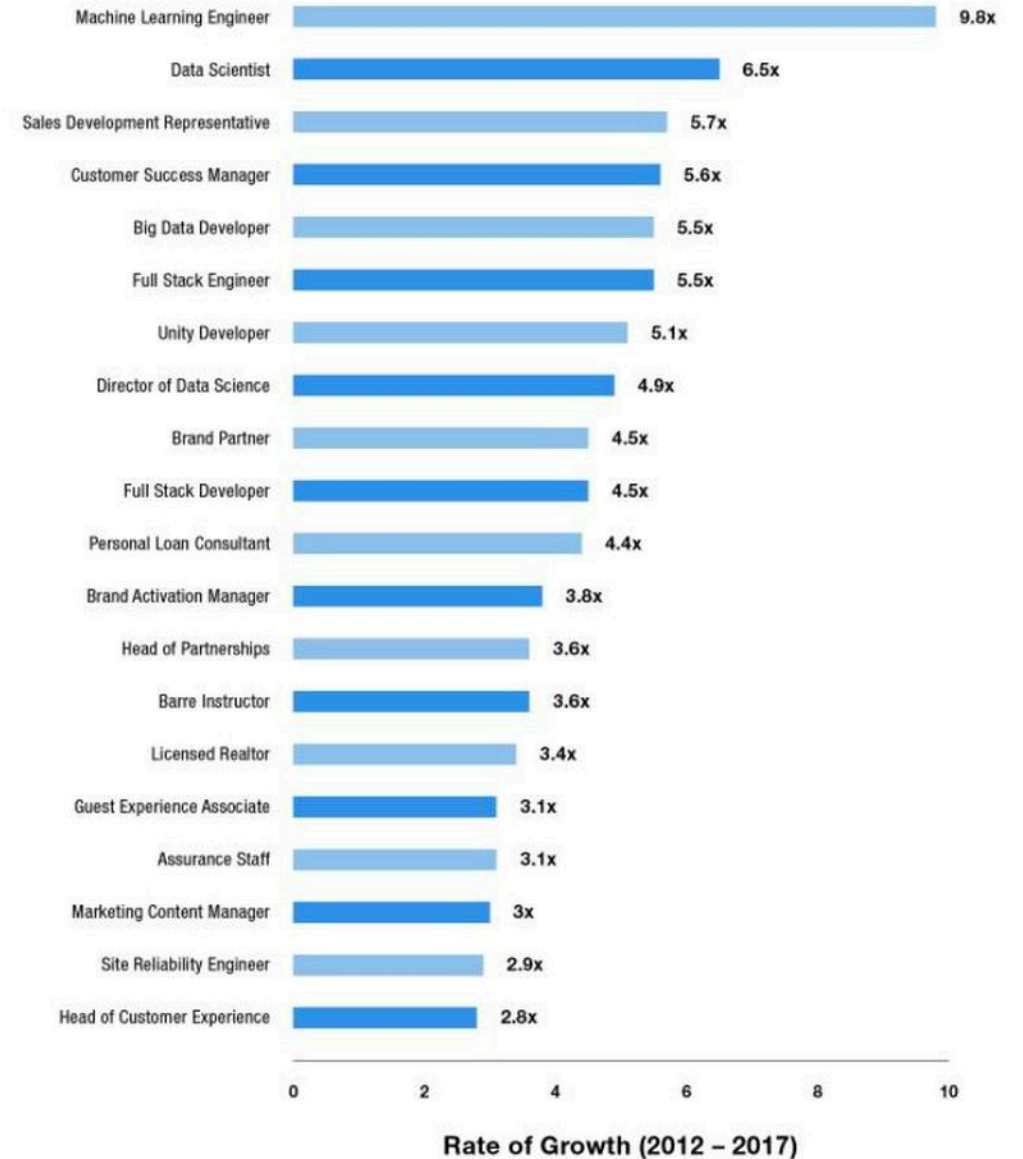


According to data

(Go see Moneyball)

# Data Science as a Profession





Top 20 Emerging Jobs — LinkedIn Economic Graph

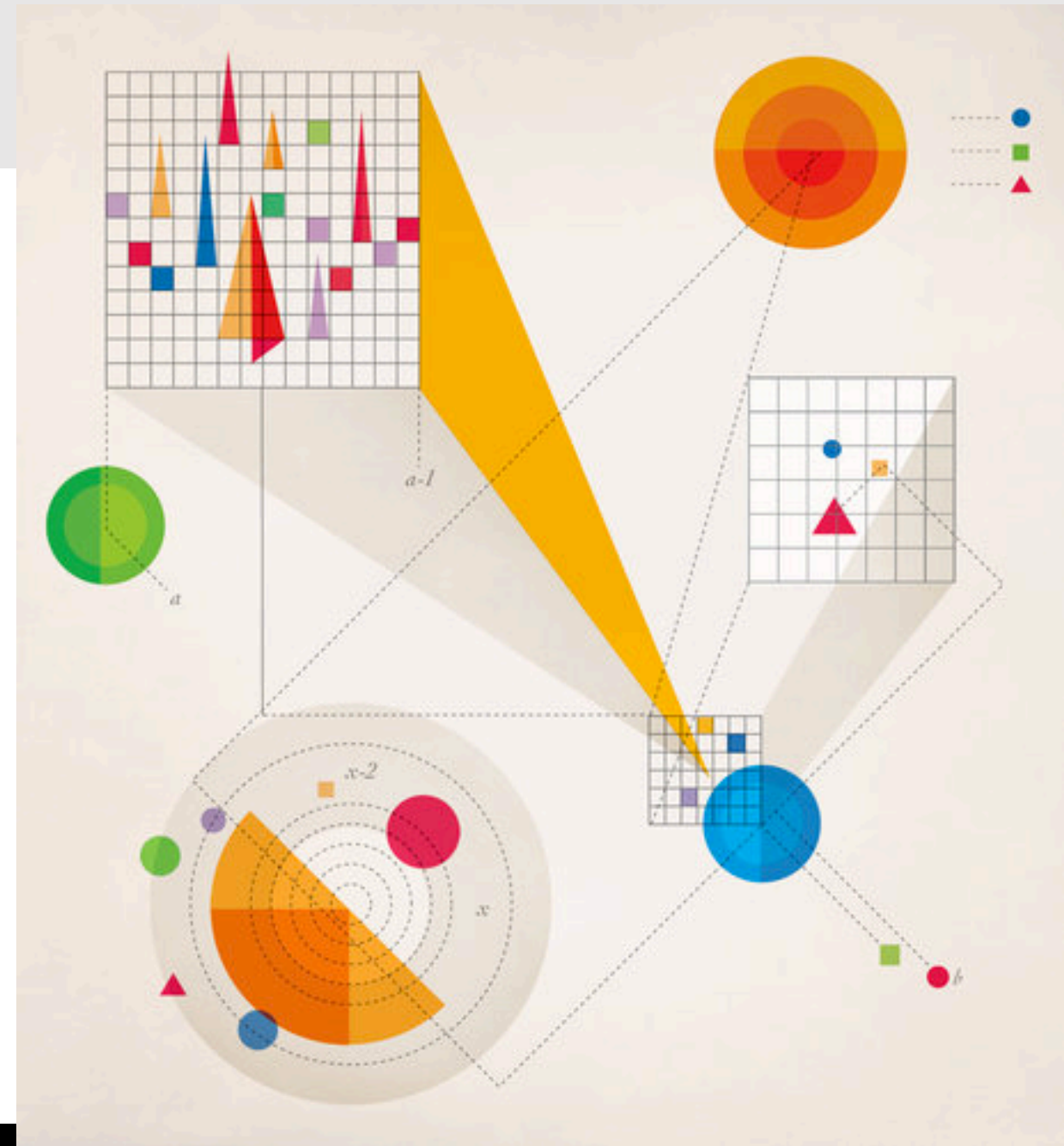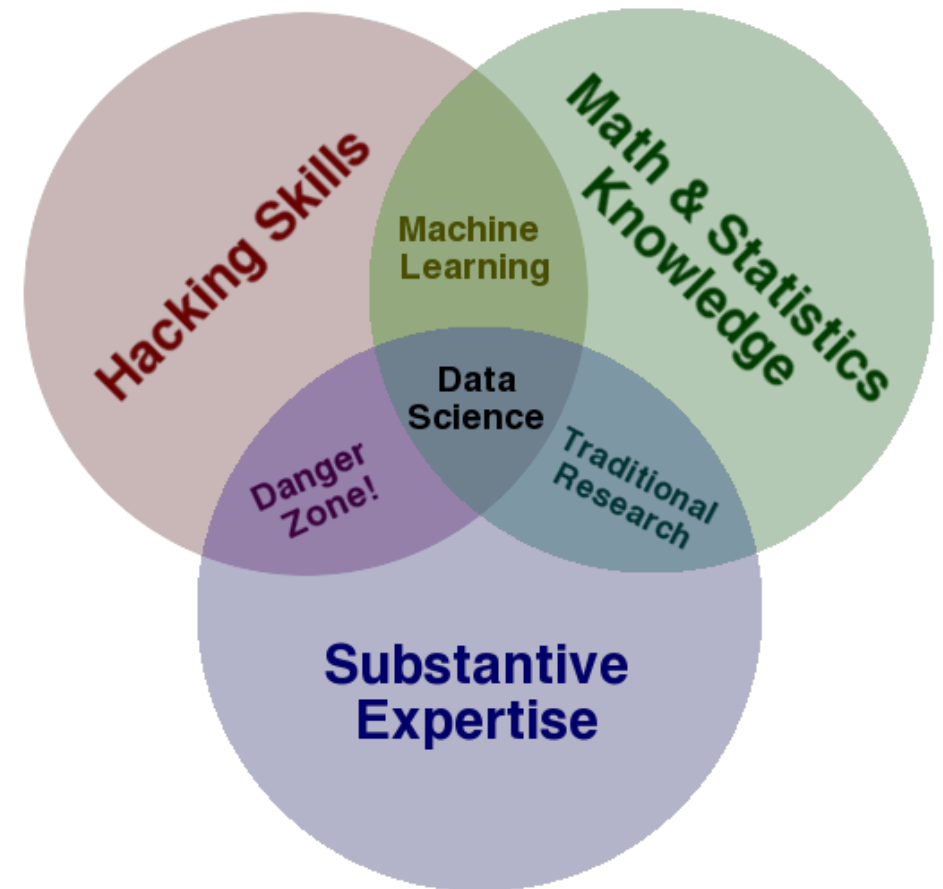| Job | Rate of Growth (2012 – 2017) |
|---|---|
| Machine Learning Engineer | 9.8x |
| Data Scientist | 6.5x |
| Sales Development Representative | 5.7x |
| Customer Success Manager | 5.6x |
| Big Data Developer | 5.5x |
| Full Stack Engineer | 5.5x |
| Unity Developer | 5.1x |
| Director of Data Science | 4.9x |
| Brand Partner | 4.5x |
| Full Stack Developer | 4.5x |
| Personal Loan Consultant | 4.4x |
| Brand Activation Manager | 3.8x |
| Head of Partnerships | 3.6x |
| Barre Instructor | 3.6x |
| Licensed Realtor | 3.4x |
| Guest Experience Associate | 3.1x |
| Assurance Staff | 3.1x |
| Marketing Content Manager | 3x |
| Site Reliability Engineer | 2.9x |
| Head of Customer Experience | 2.8x |

# Data-Literate



[McKinsey Global Institute](#) projected that the United States needs 140,000 to 190,000 more workers with "deep analytical" expertise and 1.5 million more data-literate managers, whether retrained or hired.

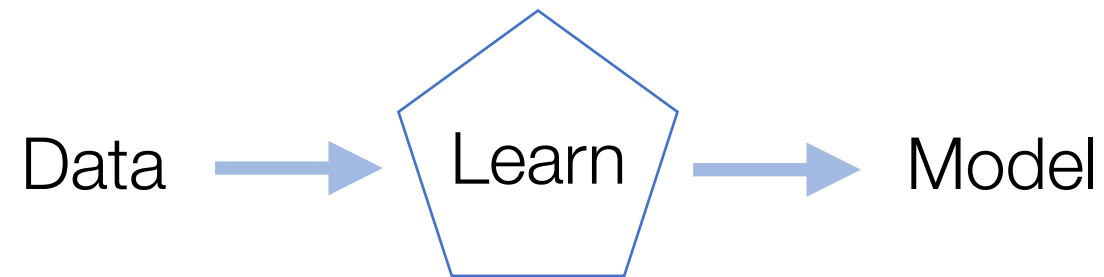http://www.nytimes.com/2012/02/12/sunday-review/big-datas-impact-in-the-world.html

# What is Data Science?

Data science is a professional approach to apply data engineering, statistics, and machine learning to solve problems in a scientific way
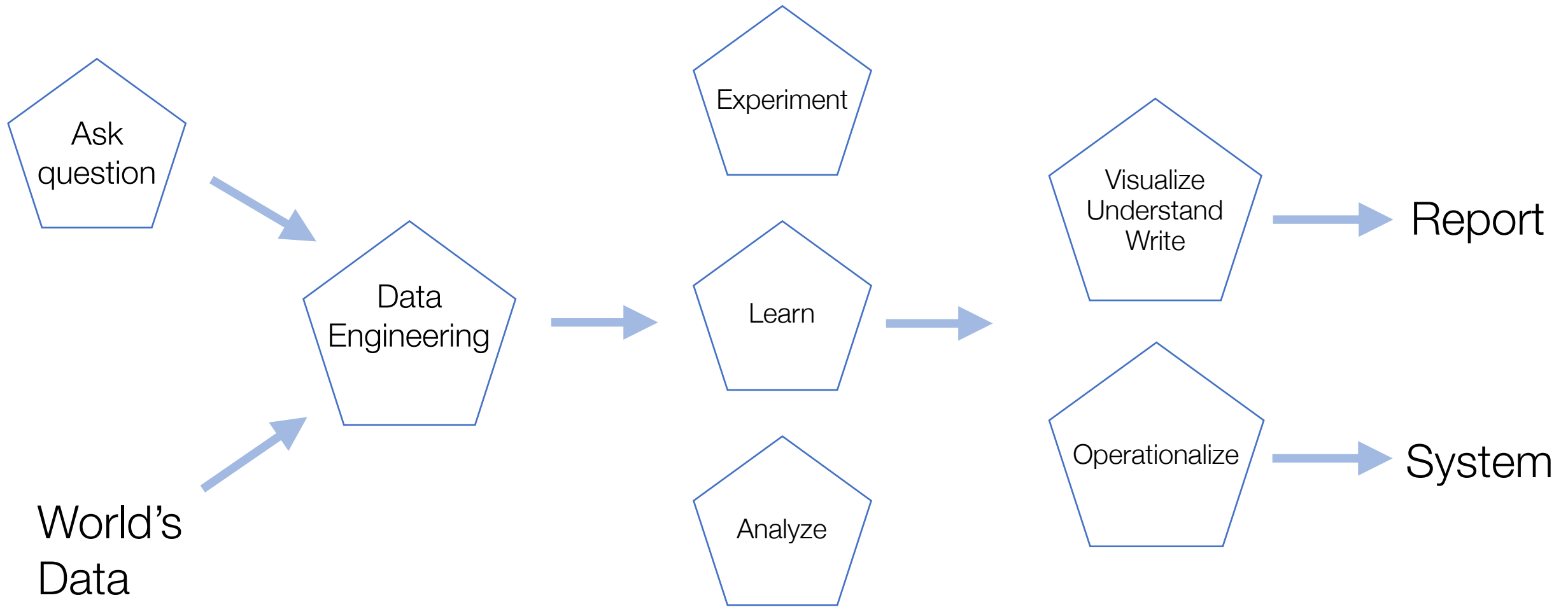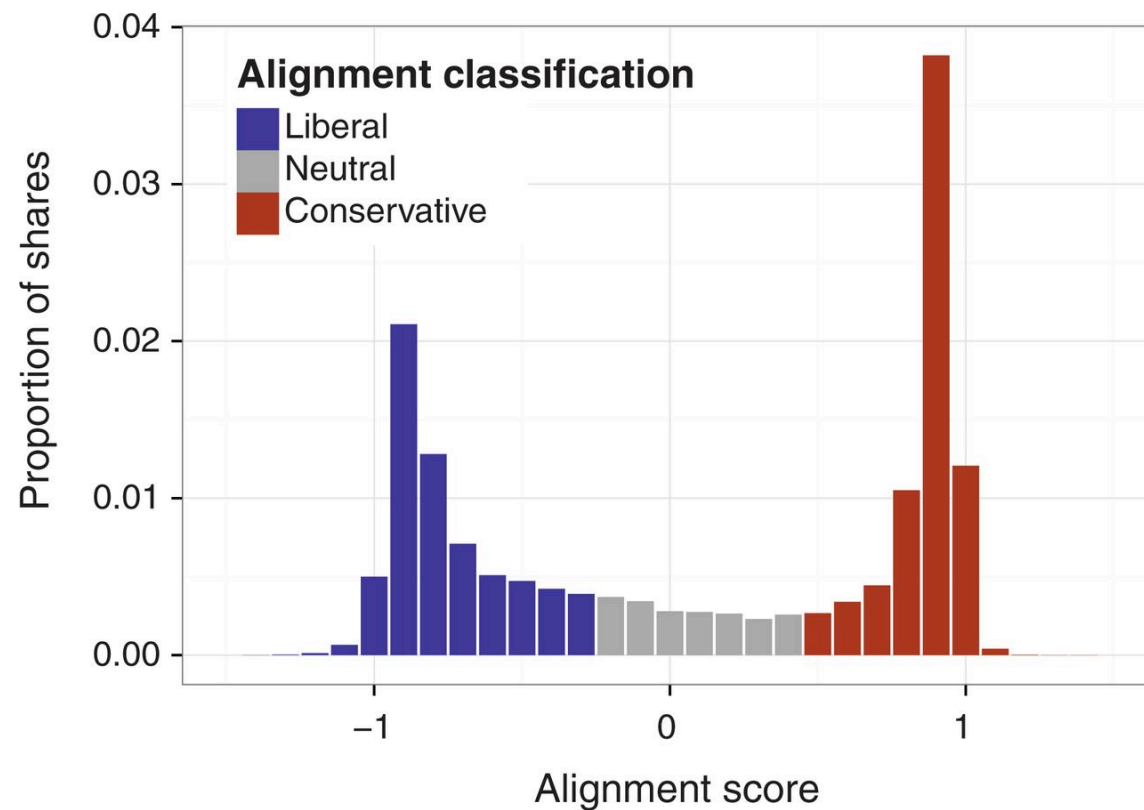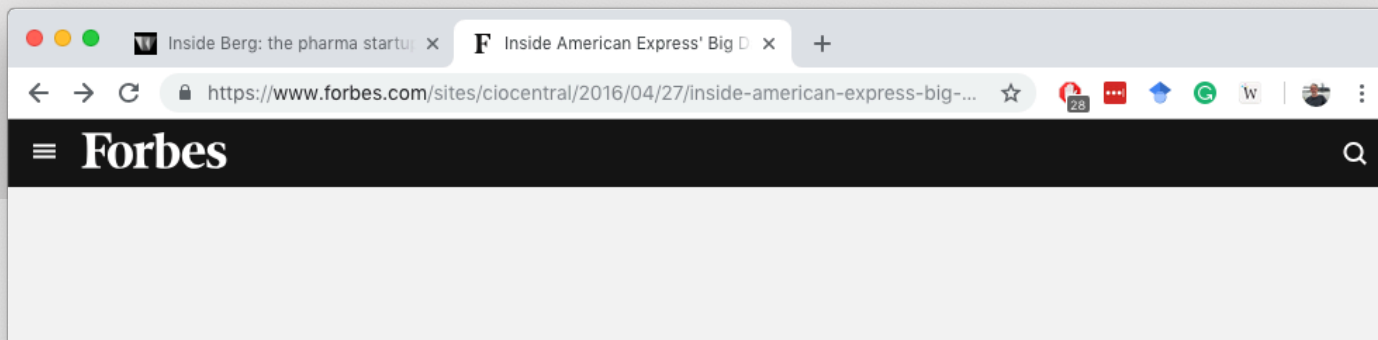
Hacking Skills

Math & Statistics Knowledge

Machine Learning

Data Science

Danger Zone!

Traditional Research

Substantive Expertise

http://drewconway.com/zia/2013/3/26/the-data-science-venn-diagram

# Buzz word hell

- Data science is a heavily criticized concept
- It is hard to distinguish it from science
- And from any type of data-intensive transaction

# The Machine Learning Model

Data →  Learn  → Model

# The Data Science Model

# Science

**Browser content:**

**Forbes**

21,700 views | Apr 27, 2016, 02:54pm

# Inside American Express' Big Data Journey

**CIO Central Guest** Contributor
**CIO Network** Contributor Group ⓘ

POST WRITTEN BY

**Randy Bean**

Randy Bean is CEO and managing partner of consultancy NewVantage Partners. You can follow him at **@RandyBeanNVP**.
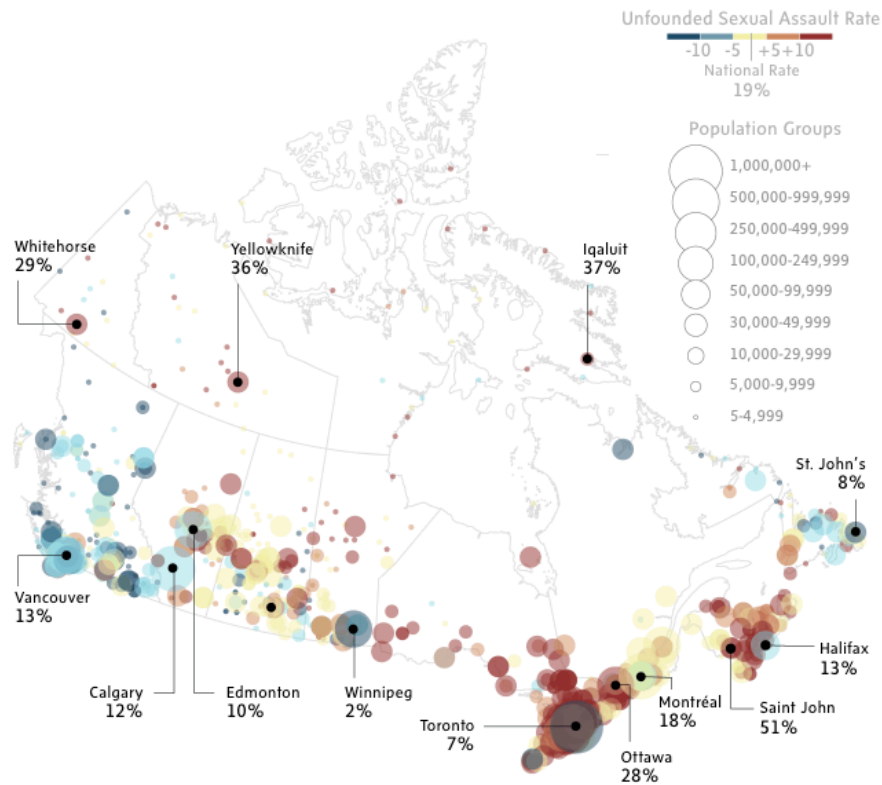
- Risk 2020 -- American Express conceptualizes how the economy and marketplace might evolve in the coming years and what are the most important risk capabilities to maintain to proactively address the weakness in the economy, a steady move towards mobile computing, cloud, artificial intelligence and deep-learning.
- Cornerstone -- This is a global, big data ecosystem is where data is organized in one place with shared global capabilities, to democratize its use across functions and geographies, recognizing that the very essence of innovation must happen at the company's DNA rather than exclusively from the top.
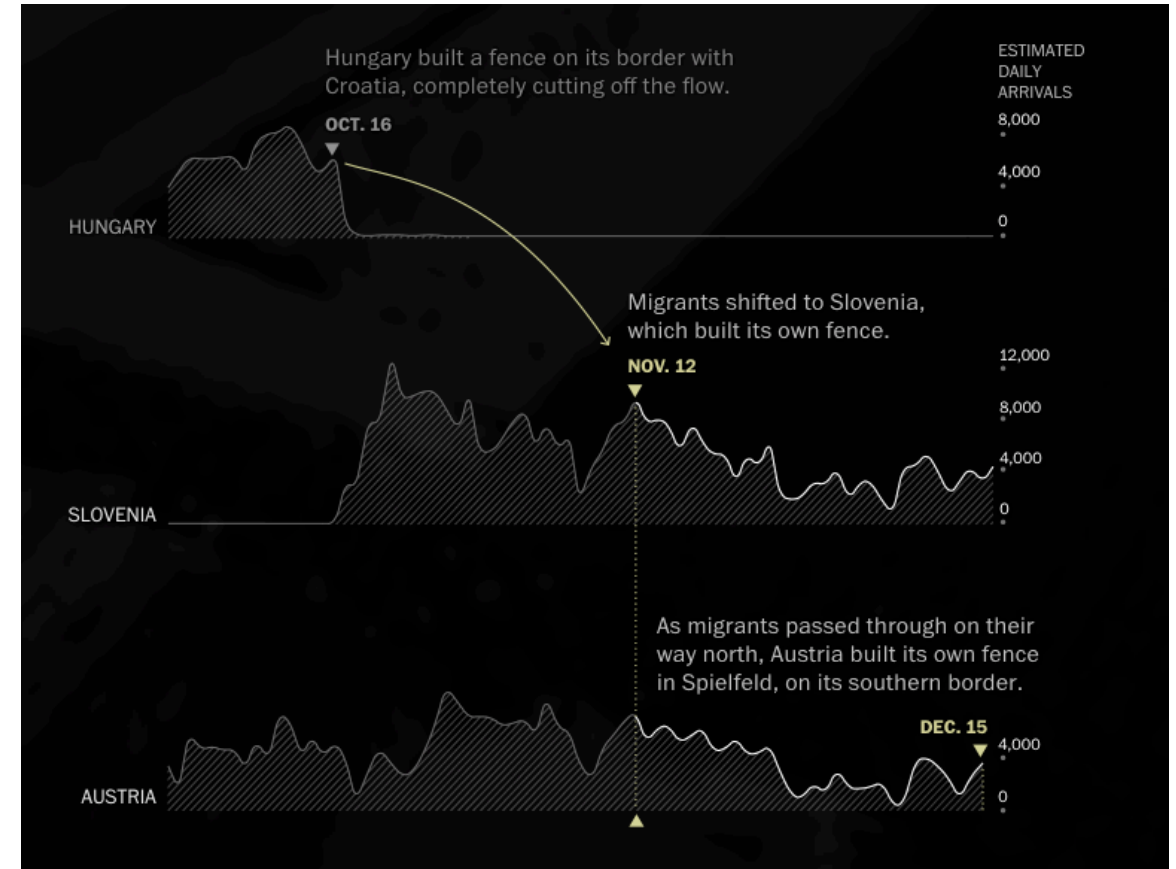
# Data Pharmaceuticals

For example, researchers at biotechnology company Berg, near Boston, Massachusetts, have developed a model to identify previously unknown cancer mechanisms using tests on more than 1,000 cancerous and healthy human cell samples. They modelled diseased human cells by varying the levels of sugar and oxygen the cells were exposed to, and then tracked their lipid, metabolite, enzyme and protein profiles. The group uses its AI platform to generate and analyse immense amounts of biological and outcomes data from patients to highlight key differences between diseased and healthy cells.

# Journalism

# Sports



https://fivethirtyeight.com/features/lionel-messi-is-impossible/



https://www.janetzko.eu/project/soccer/

# Politics



MIT Technology Review

## Intelligent Machines

## How Obama's Team Used Big Data to Rally Voters

How President Obama's campaign used big data to rally individual voters.

by Sasha Issenberg    December 19, 2012



## Cambridge Analytica, the shady data firm that might be a key Trump-Russia link, explained

Why House investigators think this company might have gamed Facebook and helped Russia spread fake news.

By Sean Illing | @seanilling | sean.illing@vox.com | Updated Apr 4, 2018, 3:41pm EDT

Photo by Bryan Bedder/Getty Images for Concordia Summit

Part of  **The Cambridge Analytica Facebook scandal**

# Summary

- Data science overwhelms science, business, and civics
- The main challenges are not technical:
  - Asking good research questions
  - Applying the right tools
  - Creating data pipelines
  - Telling a story

# The Data Science Capabilities

1. Understand the data science process

2. Model problems and answer them with <u>real</u> data

3. Control the standard "toolbox" of data science methods

4. Analyze the quality of data science results

5. Know how to report, visualize, and discuss findings

6. Introduced to the societal challenges of data science

https://hbr.org/2018/08/what-data-scientists-really-do-according-to-35-data-scientists?referral=03758&cm_vc=rr_item_page.top_right

# The Data Science Process



**Data**

# Data Engineering

ETL (Extract, Transform, and Load) is the process in which data is integrated and transferred from the operating systems to the data warehouse.



Data Staging Area

# Big Data Storage and Processing

- How to manage massive amounts of data in a way which is optimized for analysis

- Learning general data warehousing models

- Post-rational technologies: based on distributed file systems and processing:
  - Hadoop
  - Hive
  - Spark

# Experiments

- Introduction to experiment design
- Parametric and non-parametric data modeling
- Statistical tests
- Running online experiments

# The Interface with machine Learning

Understanding the interfaces with machine learning:

- Deep Sequence Learning
- Exploratory Data Analysis
- Reinforcement Learning

# The Quality of Data Science

- How to evaluate the quality of data science models?
- Identifying bias
- Simulation
- Impact assessment
- Model explainability



Multiple ROC Curves

# Reporting



- Visualization methods
  - What makes a good data visualization?
- Reporting principles and communicating data
- Operationalizing data

John Snow's map of the 1854 Broad Street cholera epidemic

# Ethics of Data Science

- Legal and ethical boundaries of data science

- Privacy

- Fairness



Bernard Parker, left, was rated high risk; Dylan Fugett was rated low risk. (Josh Ritchie for ProPublica)

**Machine Bias**

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica
May 23, 2016

# Summary

- Basic capabilities

- The course:

  - Website: https://eranto.github.io/cs5304-spring2019/

  - Slack: wild-data-science.slack.com

- The essence of the profession

# What is the "science" in Data Science?

- Data science is more than an engineering practice
- It is a professional approach that strives to embed scientific principles in data tasks
- It includes:
  - Applying a scientific method
  - Adhering (to some extent) to scientific ethical code
  - And to its culture

# The Data Science Process



Ask a question → Do Background Research → Construct a Hypothesis → Test it → Analyze Results → Communicate Findings

Do Background Research → Do Exploratory Research → Construct a Hypothesis

# Formulate a question

- Research questions should be:
  - Crunchy (either true or false)
  - Asking a question about something that can be observed: How, What, When, Who, Which, Why, or Where?
- Background research should make sure the questions should reflect the state of the art
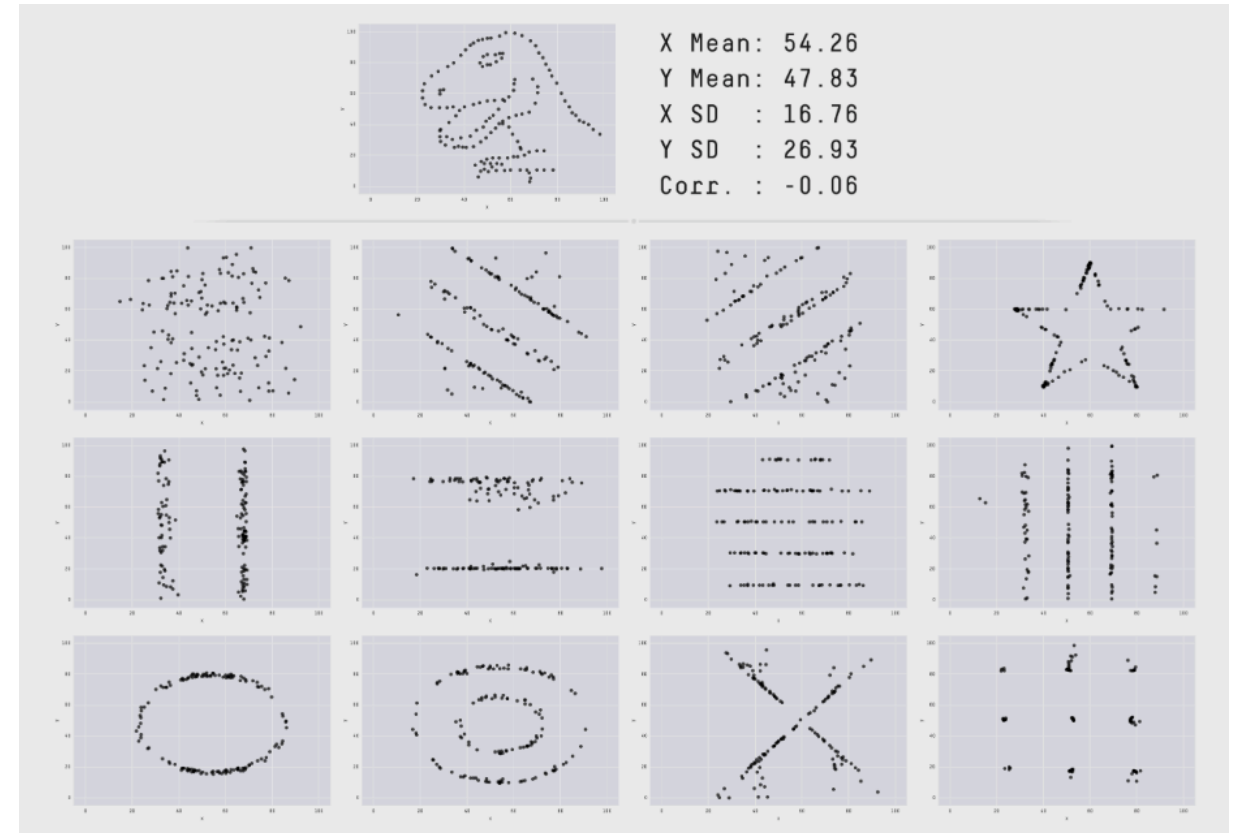
Ask a question

# Hypotheses Making

- The scientific method asks for a clearly defined hypothesis:
  - an educated guess about how things work
- An exploratory data analysis can teach us about the data, but it is not enough
- We need to show that the prediction is accurate and thus the hypothesis is supported or not

Construct a Hypothesis

# Levels of Modeling

- Classification and class probability estimation

- Regression ("value estimation")

- Similarity matching

- Clustering

- Association discovery

- Profiling

- Data reduction

- Casual modeling



X Mean: 54.26
Y Mean: 47.83
X SD  : 16.76
Y SD  : 26.93
Corr. : -0.06

https://www.autodeskresearch.com/publications/samestats

# Analyzing Results

- Be ready to fail

  - Science is about taking some risks

- Analysis should lead to something bigger than just the current problem

  - In the academia, to the construction of theory

  - In practice, to the construction of generalizable business practices

Analyze Results

# Communicating Findings

- Description of the hypotheses (so readers will know what had failed)
- Review of the state of the art
- Comprehensive description of the method
  - The standard is **reproducibility**
- Explanations of measures
- The actual findings, in a way that is both truthful and appropriate to the audience
- A discussion of the meaning of the findings

Communicate Findings

# Thinking about the writing

- What is the problem?

- Why is it interesting and important?

- Why is it hard? (E.g., why do naive approaches fail?)

- Why hasn't it been solved before? (Or, what's wrong with previous proposed solutions? How does mine differ?)

- What are the key components of my approach and results?What are the limitations?

https://cs.stanford.edu/people/widom/

# The Data Science Process