

DATA SCIENCE

CLASS 1: INTRO TO DATA SCIENCE

Rob Hall

DAT SF 19 // November 30, 2015

AGENDA

I. WHAT IS A DATA SCIENTIST?

II. WHAT IS DATA SCIENCE?

III. THE DATA SCIENCE WORKFLOW

LAB:

IV. WORKING AT THE UNIX COMMAND LINE

I. WHAT IS A DATA SCIENTIST?

WHAT IS A DATA SCIENTIST?



Zvi

@nivertech



Follow

"Data Scientist" is a Data Analyst who lives in California.

RETWEETS

162

LIKES

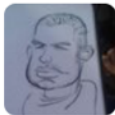
79



6:55 PM - 14 Mar 2012



WHAT IS A DATA SCIENTIST?



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.

RETWEETS

1,255

LIKES

714



9:55 AM - 3 May 2012

▲ ▲ ▲

WHAT IS A DATA SCIENTIST?



Javier Nogales

@fjnogales



Follow

Data Scientist (2/2): person who is worse at statistics than any statistician and worse at software engineering than any software engineer

RETWEET

1

LIKES

5



6:08 AM - 27 Jan 2014



WHAT MAKES A GOOD DATA SCIENTIST?



Michael E. Driscoll

@medriscoll



Following

Data scientists: better statisticians than
most programmers & better programmers
than most statisticians bit.ly/NHmRqu
[@peteskomoroch](#)



Reply



Retweet



Favorite



More



Pocket

WHAT IS YOUR DEFINITION?

“Data Scientists are people with some mix of **coding and statistical skills** who work on **making data useful** in various ways.”

Data Scientist Type A (for Analysis):

- Primarily concerned with **making sense of data** or working with it in a fairly **static** way.
- Similar to a statistician, but knows all the **practical details of working with data** that aren't taught in statistics: data cleaning, dealing with large data sets, visualization, domain knowledge, etc.

“Data Scientists are people with some mix of **coding and statistical skills** who work on **making data useful** in various ways.”

Data Scientist Type B (for Building):

- Some statistical background, but **strong coder or software engineer**.
- Primarily concerned with **using data “in production”**: building models which interact with users (by giving recommendations, for example).

Our course is focused primarily on **Type A**.

Hadley Wickham's advice for becoming a data scientist:

Statistical knowledge

“I think you need some knowledge of specific statistical/machine learning techniques, but a deep theoretical understanding is not that important. You need to understand the strengths and weaknesses of each technique... The vast majority of data science problems can be solved by a creative assembly of off-the-shelf techniques, and don't require new theory.”

Hadley Wickham's advice for becoming a data scientist:

Programming skills

“You need to be fluent with either R or Python. There are other options, but none of them have the community that R and Python have, which means you'll need to spend a lot of time reinventing tools that already exist elsewhere.”

Hadley Wickham's advice for becoming a data scientist:

Domain knowledge

“...A data scientist should be able to contribute meaningfully to any project, even if you're not intimately familiar with the specifics. I think this means you should be generally well read... and an able communicator. A good data scientist will help the real domain experts refine and frame their questions in a helpful way. Unfortunately I don't know of any good resources for learning how to ask questions.”

Chris Volinsky (Columbia & AT&T Labs) on “Data Mining vs. Statistics”

- Snark: Data Mining = Statistics + Marketing
- Statistics is known for: **well-defined hypotheses** used to learn about a **specifically chosen population** studied using **carefully collected data** providing inferences with **well-known properties**.
- Data mining isn't that careful. It is: **data-driven discovery** of **models and patterns** from **massive and observational data sets**.

II. WHAT IS DATA SCIENCE?

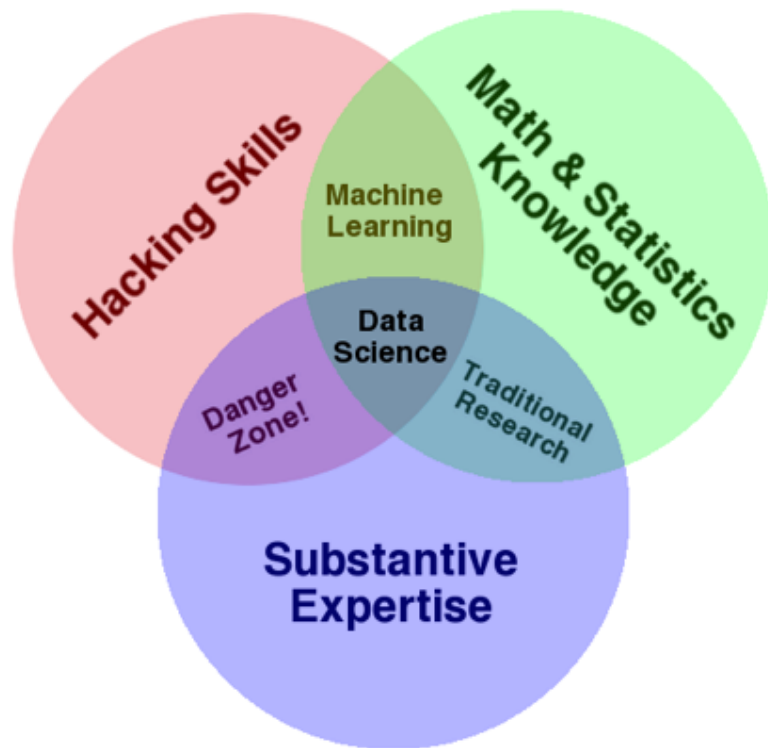
WHAT IS DATA SCIENCE?

- A set of tools and techniques used to extract useful information from data.

WHAT IS DATA SCIENCE?

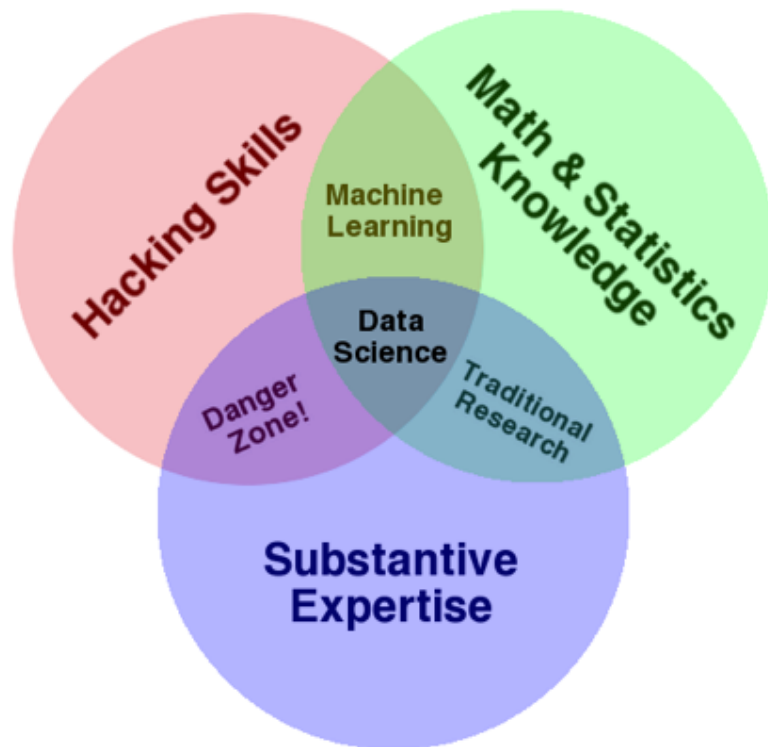
- A set of tools and techniques used to extract useful information from data.
- An interdisciplinary, problem-oriented subject.

THE QUALITIES OF A DATA SCIENTIST



source: <http://www.dataists.com/2010/09/the-data-science-venn-diagram/>

THE QUALITIES OF A DATA SCIENTIST



ONE MORE THING!

Communication skills

WHAT IS DATA SCIENCE?

- A set of tools and techniques used to extract useful information from data.
- An interdisciplinary, problem-solving oriented subject.
- The application of scientific techniques to practical problems.

WHAT IS DATA SCIENCE?

- A set of tools and techniques used to extract useful information from data.
- An interdisciplinary, problem-solving oriented subject.
- The application of scientific techniques to practical problems.
- A rapidly growing field.

HB.R.ORG

Harvard Business Review



OCTOBER 2012
REPRINT R02100

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

22

McKinsey
estimates
140,000-190,000
shortage by 2018

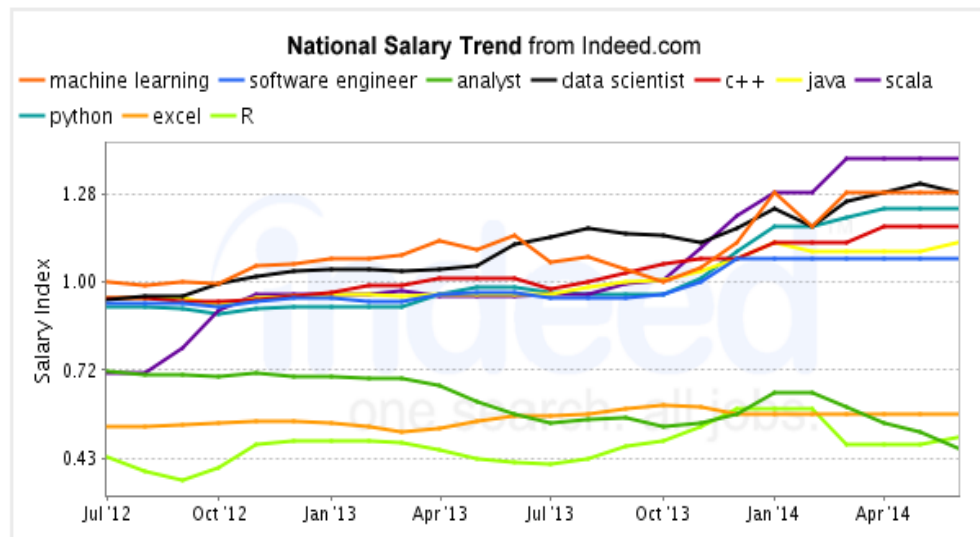
I keep saying the sexy job in the next ten years will be statisticians. People think I'm joking, but who would've guessed that computer engineers would've been the sexy job of the 1990s?

Hal Varian, Chief Economist at Google, [The McKinsey Quarterly, January 2009](#)

machine learning in San Francisco, CA	\$152,000	
software engineer in San Francisco, CA	\$139,000	
analyst in San Francisco, CA	\$72,000	
data scientist in San Francisco, CA	\$160,000	
c++ in San Francisco, CA	\$141,000	
java in San Francisco, CA	\$139,000	
scala in San Francisco, CA	\$163,000	
python in San Francisco, CA	\$146,000	
excel in San Francisco, CA	\$75,000	
R in San Francisco, CA	\$68,000	


In USD as of Nov 17, 2014

60k 120k 180k



DATA SCIENTISTS WANTED

25



Data Scientist

Facebook is seeking a Data Scientist to be comfortable working as a Data Scientist. We have a keen interest in the study of data and the questions that help us build the future.

About this job

Job description


Facebook is seeking a Data Scientist to be comfortable working as a Data Scientist. We have a keen interest in the study of data and the questions that help us build the future.

Responsibilities

- Work closely with a product manager to understand the product vision and requirements.
- Answer product questions and deliver actionable findings to the product manager.
- Drive the collection of new data and analyze and interpret the results.
- Develop best practices for product engineering team.

Requirements

- M.S. or Ph.D. in a relevant field.
- Extensive experience with data analysis and machine learning.
- Comfort manipulating large datasets.
- A strong passion for engineering and a flexible analytic approach.
- Ability to communicate effectively with non-technical stakeholders.
- Fluency with at least one of the following: Python, R, Java, or C++.
- Familiarity with relational databases and NoSQL stores.
- Expert knowledge of an analytical tool like Hadoop, MapReduce, or Hive.
- Experience working with large datasets and building scalable data processing pipelines.



Data Scientist

EMC is seeking a Data Scientist to be comfortable working as a Data Scientist. We have a keen interest in the study of data and the questions that help us build the future.

About this job

Job description


EMC is seeking a Data Scientist to be comfortable working as a Data Scientist. We have a keen interest in the study of data and the questions that help us build the future.

Responsibilities

- Work closely with a product manager to understand the product vision and requirements.
- Answer product questions and deliver actionable findings to the product manager.
- Drive the collection of new data and analyze and interpret the results.
- Develop best practices for product engineering team.

Requirements

- M.S. or Ph.D. in a relevant field.
- Extensive experience with data analysis and machine learning.
- Comfort manipulating large datasets.
- A strong passion for engineering and a flexible analytic approach.
- Ability to communicate effectively with non-technical stakeholders.
- Fluency with at least one of the following: Python, R, Java, or C++.
- Familiarity with relational databases and NoSQL stores.
- Expert knowledge of an analytical tool like Hadoop, MapReduce, or Hive.
- Experience working with large datasets and building scalable data processing pipelines.



Data Scientist at LinkedIn

LinkedIn - Mountain View, CA

Posted 26 days ago

About this job

Job description

As a Senior Data Scientist at LinkedIn, you will develop innovative new technologies, features, and products that help connect the world's professionals to make them more productive and successful.

Description


Our team applies machine learning techniques on social data to build products & features that reach over 200M professionals on LinkedIn. We build graph and text mining systems to tackle hard problems in areas like entity resolution, search relevance, recommendation algorithms, reputation & skills assessment, and network analysis.

Along with our team of data scientists, you'll work with product managers, designers, and engineers to build data driven features and products like LinkedIn Skills, Endorsements, and InMaps.

If you enjoy working with data to build products and solve hard problems in creative ways, you will fit right in.

Requirements:

- Strong background in Machine Learning, Statistics, Information Retrieval, or Graph Analysis
- Some experience working with large datasets, preferably using tools like Hadoop, MapReduce, Pig, or Hive
- 2+ years experience developing high quality software, contributions to open source projects are a plus
- Experience programming in an object oriented language (Java, C++, etc)
- Knowledge of scripting languages like Ruby or Python, familiarity with web frameworks a plus
- Comfortable with data analysis & visualization using tools like R, Matlab, or SciPy
- Critical thinking: ability to track down complex data and engineering issues, evaluate different algorithmic approaches, and analyze data to solve problems
- Creativity: you can conceive of new data driven products, features, and technologies
- Results: you prioritize, focusing on ideas and features that will have significant, measurable impact
- Planning & estimation: ability to set and meet your own project objectives & milestones
- Ability to coordinate effectively with team members in engineering, design, and product management
- Communicate results and progress internally and externally in meetings, presentations, and tech talks
- Masters, PhD, or equivalent experience in a quantitative field (computer science, physics, mathematics, bioinformatics, etc.)



Data Scientist

Apple - Santa Clara Valley - California - US

Posted 19 days ago

About this job

Job description

Apple has a tremendous amount of data, and we have just scratched the surface in pattern detection, anomaly detection, predictive modeling, and optimization. There are many exciting problems to be discovered and solved. We encourage scientists to stay abreast of data mining research by attending conferences and working with academic faculty and students. We foster a collaborative work environment, but allow solution autonomy on projects.

The iTunes Engineering team has a proud tradition of delivering cutting-edge products in a competitive marketplace. We seek to maintain a challenging and rewarding environment where the best engineers and scientists can collaborate and produce real-world improvements in customers' online experience. Successful candidates will solve problems unique in scale and concept in the pursuit of new and original features.

Key Qualifications

- Strong working knowledge of data mining algorithms including decision trees, probability networks, association rules, clustering, regression, and neural networks.
- Familiarity with database modeling and data warehousing principles with a working knowledge of SQL.
- Familiarity with Big Data tools and techniques, including MapReduce, NoSQL stores, and unbounded stream processing.
- Creativity to go beyond current tools to deliver best solution to the problem
- Strong programming skills in Java, Python, or similar language
- Excellent interpersonal, written, and verbal communication skills
- Ability and comfort working independently and making key decisions on projects

Description

We are seeking an outstanding data mining scientist who is interested in designing, developing, and fielding data mining solutions that have direct and measurable impact to Apple. This person will work within and across teams to help identify viable data mining opportunities and then implement end to end analytical solutions. The role requires both a broad knowledge of existing data mining algorithms and creativity to invent and customize when necessary.

Education

Ph.D. in Data Mining, Machine Learning, Statistics, Operations Research or related field

M.S. in related field with 5 years experience applying data mining techniques to real business problems.

WHO USES DATA SCIENCE?



WHAT MAKES A GOOD DATA SCIENTIST?

- Statistical and machine learning knowledge
- Engineering experience
- Curiosity
- Product sense
- Storytelling
- Cleverness

WHO ARE DATA SCIENTISTS?

Figure 8. Data Scientists by Area of Study

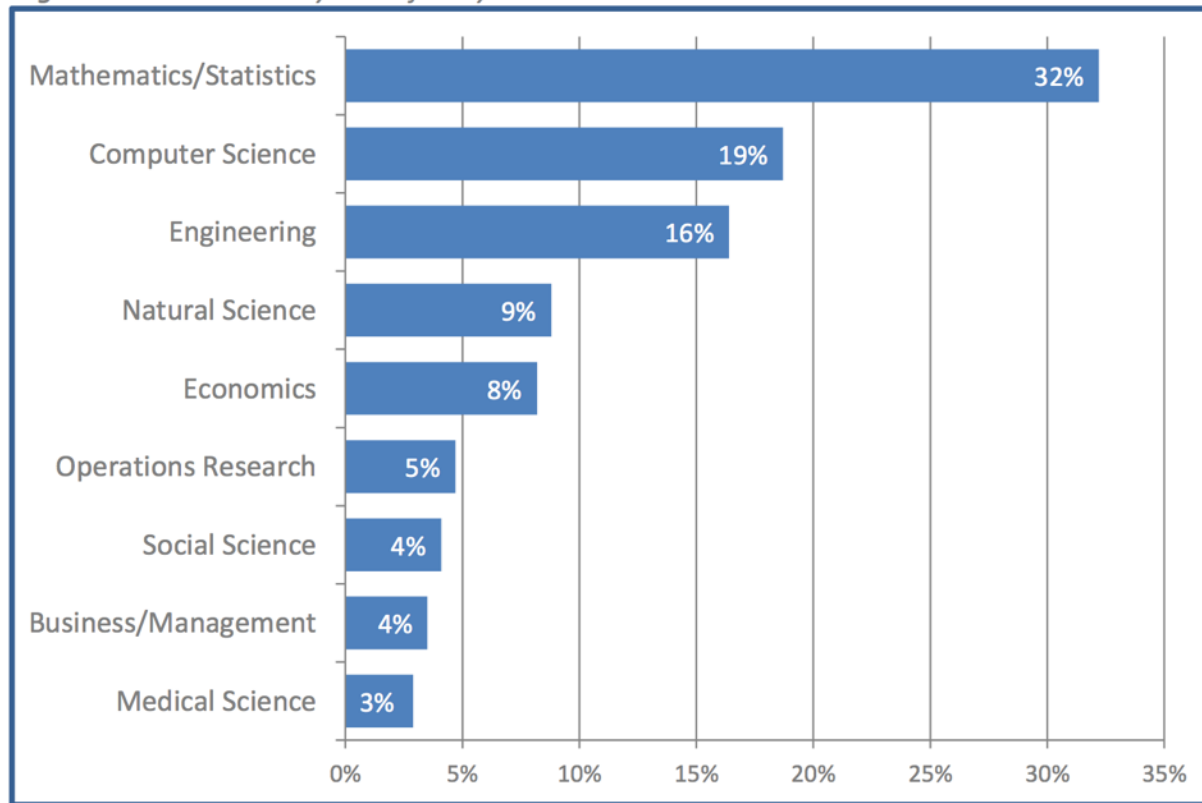
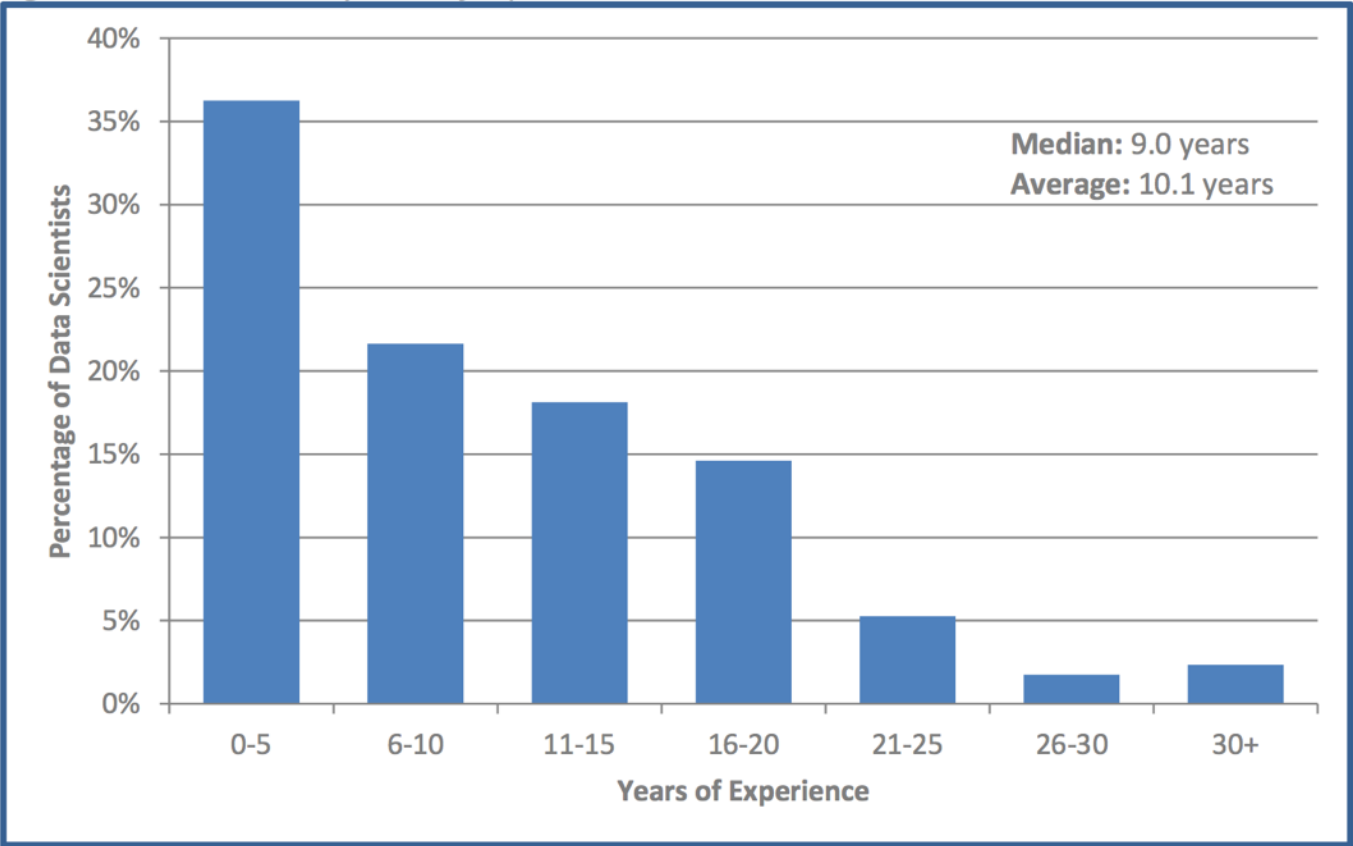
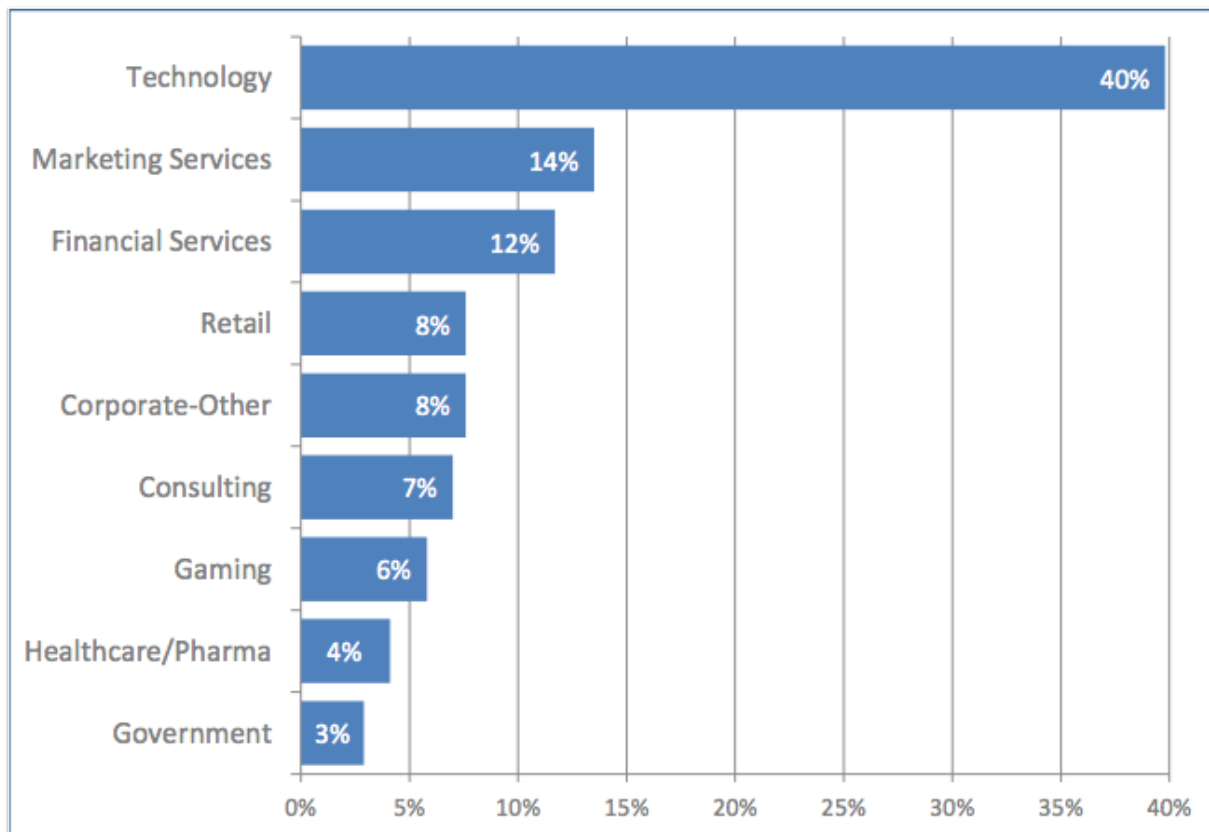
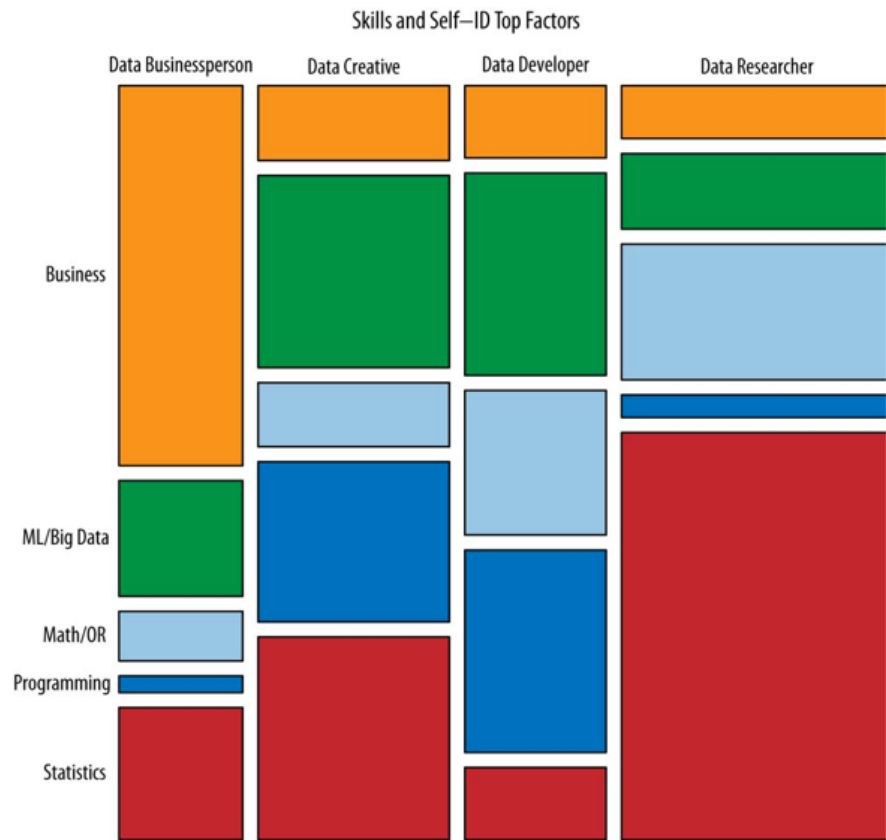


Figure 4. Data Scientists by Years of Experience



WHO ARE DATA SCIENTISTS?



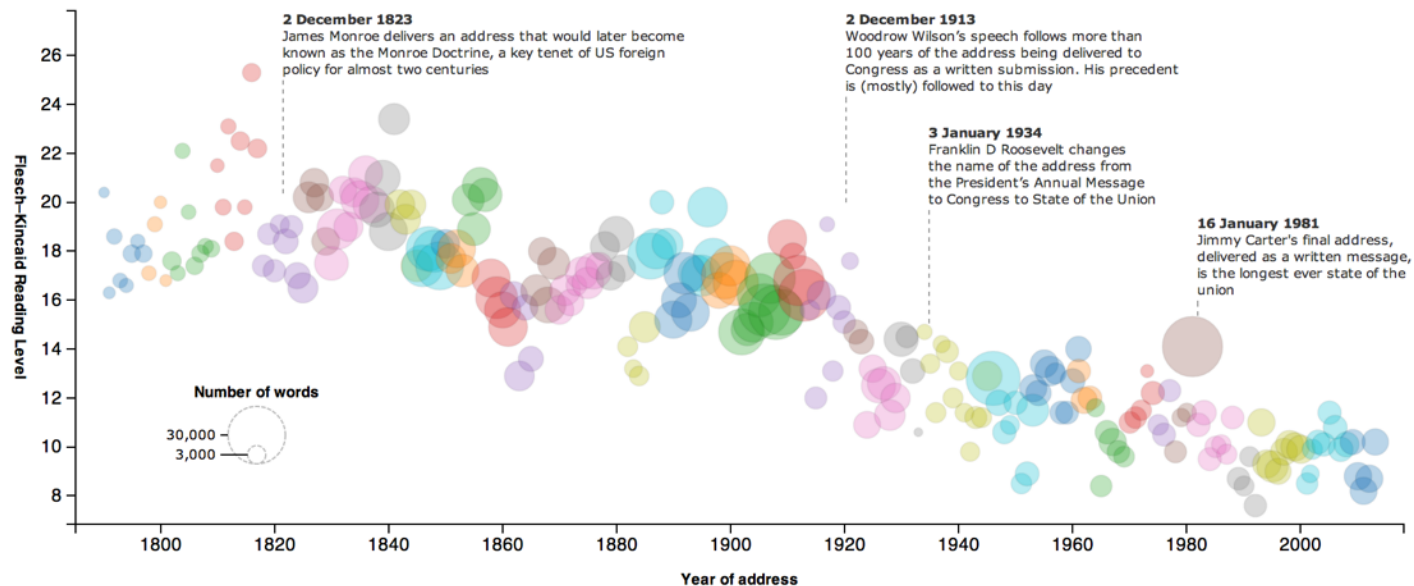


WHO USES DATA SCIENCE?

The state of our union is ... dumber:

How the linguistic standard of the presidential address has declined

Using the [Flesch-Kincaid readability test](#) the Guardian has tracked the reading level of every state of the union



Music + Data:
<http://bit.ly/echonest>

III. THE DATA SCIENCE WORKFLOW

THE DATA SCIENCE WORKFLOW

Dataists (Hilary Mason & friends)

- 1. Obtain
- 2. Scrub
- 3. Explore
- 4. Model
- 5. Interpret

Dataists (Hilary Mason & friends)

- 1. Obtain - pointing and clicking does not scale (APIs, Python, shell scripting)
- 2. Scrub - “Scrubbing data is the least sexy part of the analysis process, but often one that yields the greatest benefits” (Python, sed, awk, grep)
- 3. Explore - look at the data (visualizing, clustering, dimensionality reduction)
- 4. Model - “All models are wrong, but some are useful” / models are built to predict and interpret
- 5. Interpret - “The purpose of computing is insight, not numbers”

Jeff Hammerbacher (Facebook, Cloudera)

- 1. Identify problem
- 2. Instrument data sources
- 3. Collect data
- 4. Prepare data (integrate, transform, clean, impute, filter, aggregate)
- 5. Build model
- 6. Evaluate model
- 7. Communicate results

THE DATA SCIENCE WORKFLOW

Ted Johnson

- 1. Assemble an accurate and relevant data set
- 2. Choose the appropriate algorithm

THE DATA SCIENCE WORKFLOW

Ben Fry

- 1. Acquire
- 2. Parse
- 3. Filter
- 4. Mine
- 5. Represent
- 6. Refine
- 7. Interact

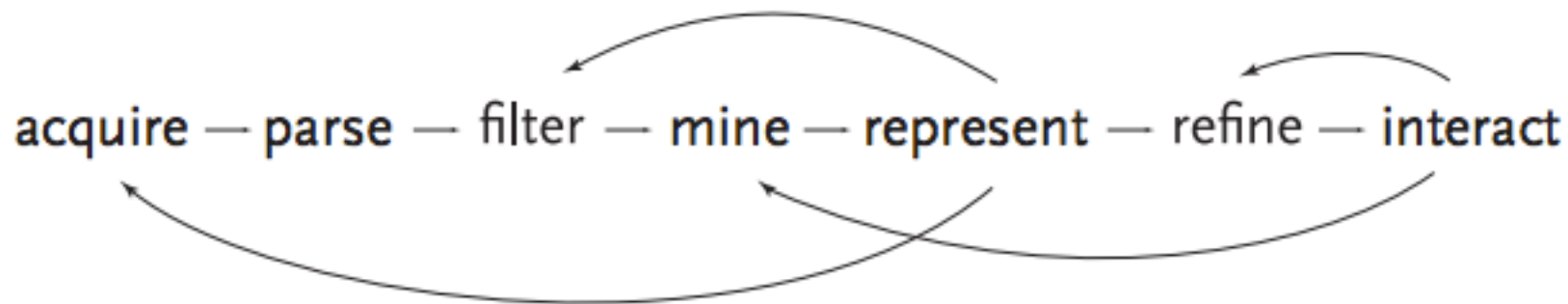
Ben Fry

- 1. Acquire - the matter of obtaining the data
- 2. Parse - providing some structure around what the data means
- 3. Filter - removing all but the data of interest
- 4. Mine - the application of methods from statistics or data mining, as a way to discern patterns or place the data in mathematical context
- 5. Represent - determination of a simple representation (e.g. graphing)
- 6. Refine - improvements to the basic representation to make it clearer and more visually engaging
- 7. Interact - the addition of methods for manipulating the data or controlling which features are visible

THE DATA SCIENCE WORKFLOW



THE DATA SCIENCE WORKFLOW



NOTE

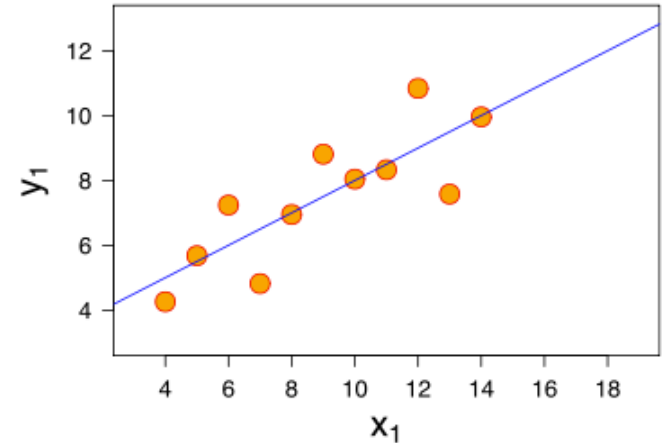
This diagram illustrates the *iterative* nature of problem solving

III. VISUALIZATIONS AS A MEDIUM

EXERCISE – WHY VISUALIZE DATA?

Consider the following dataset:

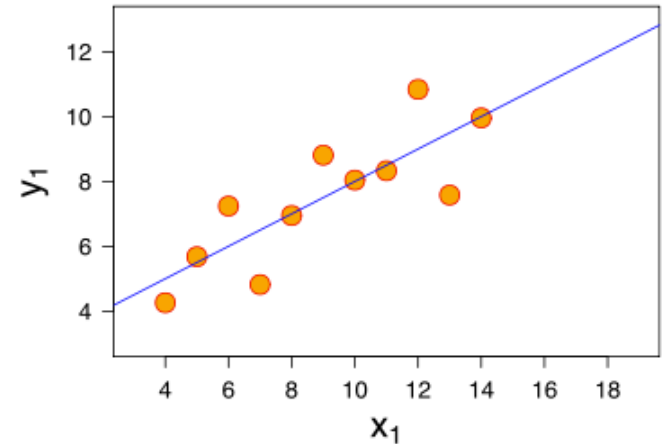
- *eleven (x, y) points*



EXERCISE – WHY VISUALIZE DATA?

Consider the following dataset:

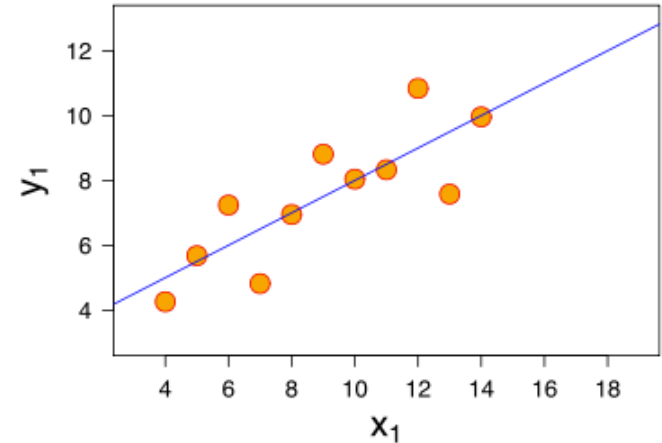
- *eleven (x, y) points*
- *mean of $x = 9$, mean of $y = 7.5$*



EXERCISE – WHY VISUALIZE DATA?

Consider the following dataset:

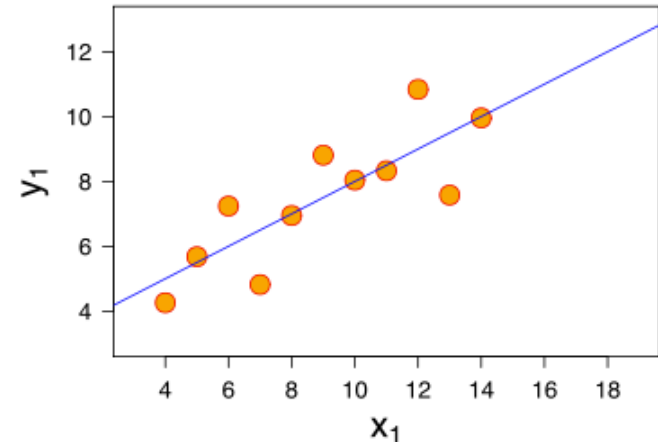
- eleven (x, y) points*
- mean of $x = 9$, mean of $y = 7.5$*
- variance of $x = 11$, variance of $y = 4.1$*



EXERCISE – WHY VISUALIZE DATA?

Consider the following dataset:

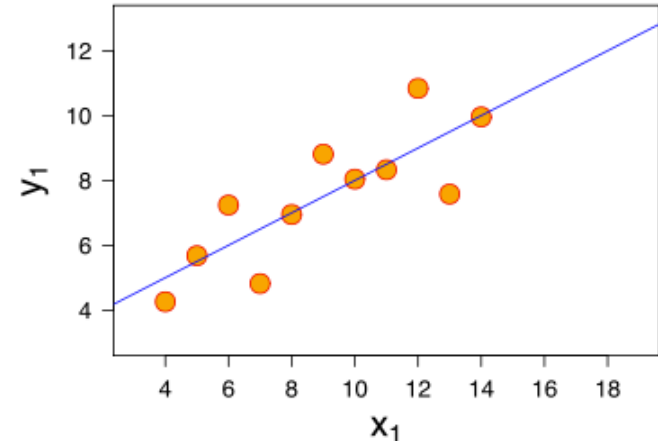
- eleven (x, y) points*
- mean of $x = 9$, mean of $y = 7.5$*
- variance of $x = 11$, variance of $y = 4.1$*
- correlation of x and $y = 0.8$*



EXERCISE – WHY VISUALIZE DATA?

Consider the following dataset:

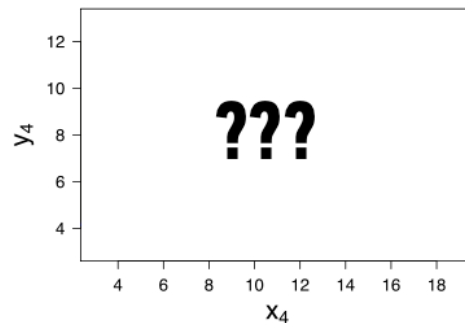
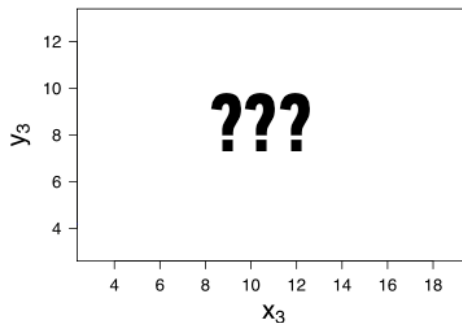
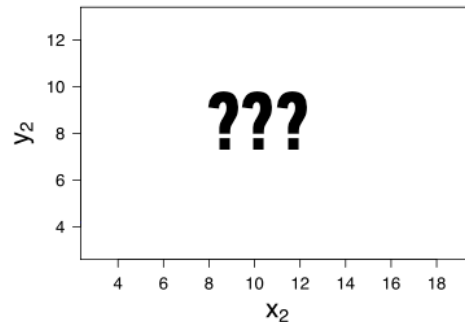
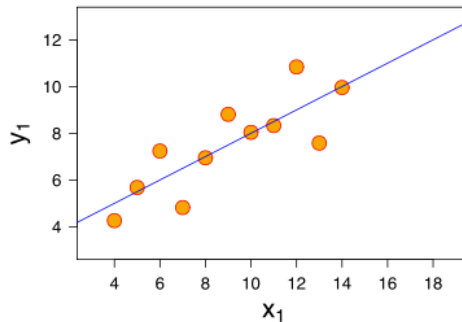
- *eleven (x, y) points*
- *mean of $x = 9$, mean of $y = 7.5$*
- *variance of $x = 11$, variance of $y = 4.1$*
- *correlation of $x, y = 0.8$*
- *line of best fit: $y = 3.00 + 0.500x$*



EXERCISE – WHY VISUALIZE DATA?

*Now, suppose I give you
three more datasets
with exactly the same
characteristics...*

*Q: how similar are these
datasets?*

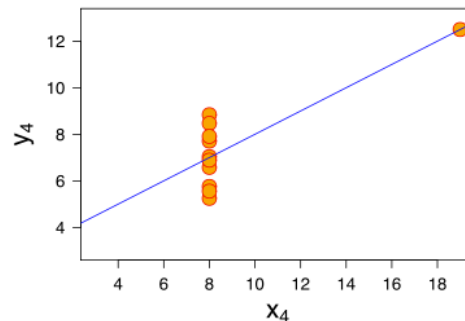
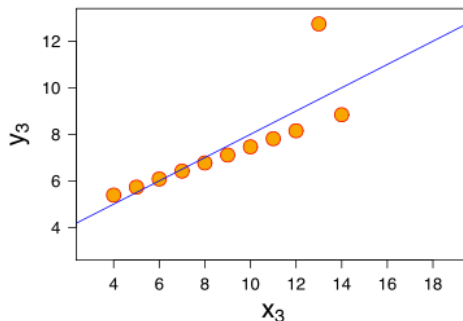
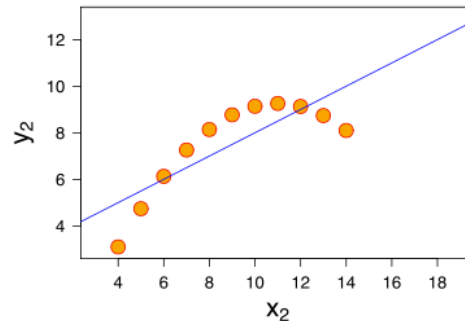
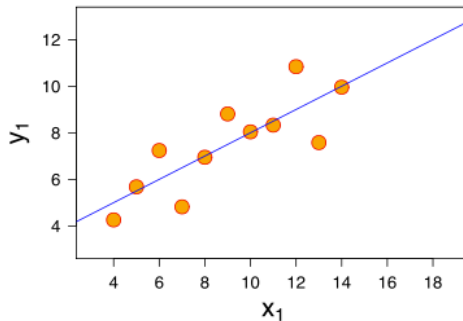


EXERCISE – WHY VISUALIZE DATA?

*Now, suppose I give you
three more datasets
with exactly the same
characteristics.*

*Q: how similar are these
datasets?*

A: not very!



EXERCISE – WHY VISUALIZE DATA?

Look at your data!

V. EXERCISE

DISCUSSION