# 菩萨相谈

菩

# Data Science

Brian Sletten

@bsletten
09/29/2014

# Speaker Qualifications

- Specialize in next-generation technologies
- Author of "Resource-Oriented Architecture Patterns for Webs of Data"
- Speaks internationally about REST, Semantic Web, Security, Visualization, Architecture
- Worked in Defense, Finance, Retail, Hospitality, Video Game, Health Care and Publishing Industries
- One of Top 100 Semantic Web People

# Agenda

- Introduction
- Data Science Techniques
- Programming
- Visualization
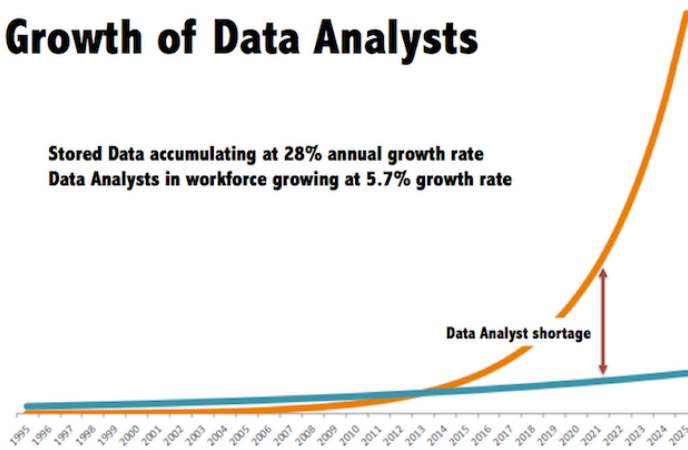- Machine Learning
- Data Mining
- Big Data
- Linked Data

菩

# Introduction

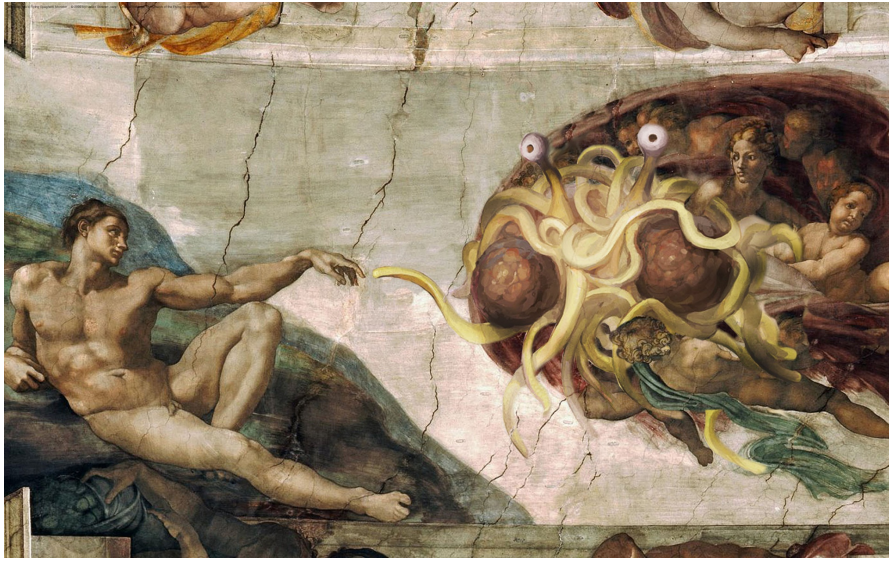**Growth of Data vs. Growth of Data Analysts**

Stored Data accumulating at 28% annual growth rate
Data Analysts in workforce growing at 5.7% growth rate

Data Analyst shortage

1995 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020 2021 2022 2023 2024 2025

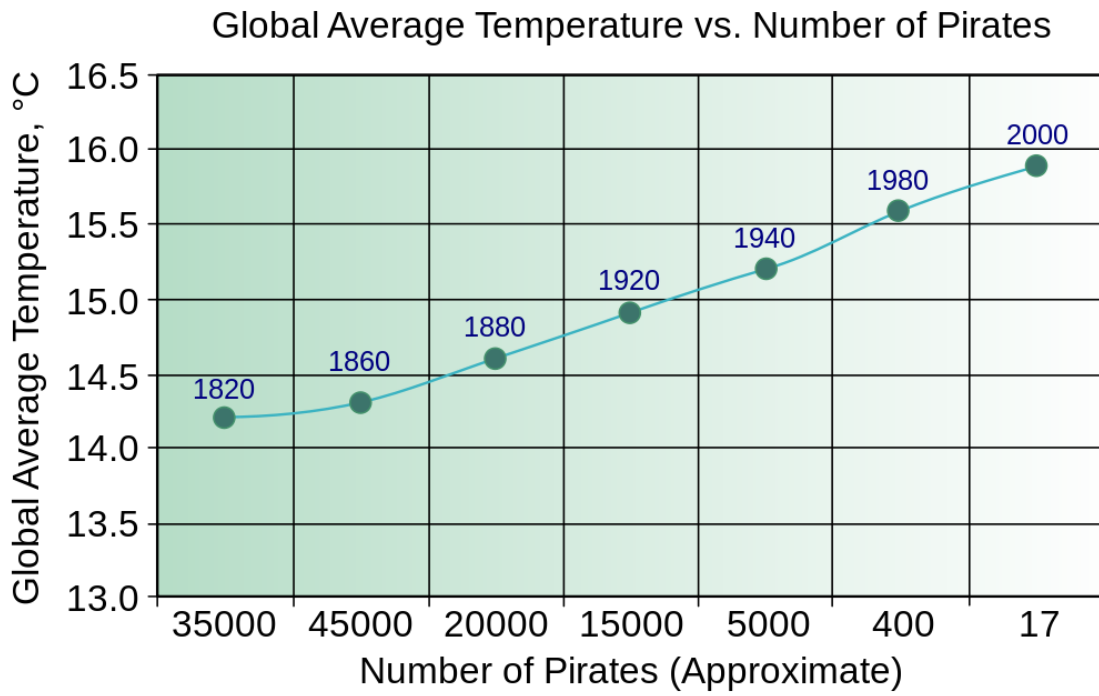http://www.delphianalytics.net/wp-content/uploads/2013/04/GrowthOfDataVsDataAnalysts.png

"We're witnessing the beginning of a massive, culturally saturated feedback loop

where our behavior changes the product and the product changes our behavior.

Technology makes this possible: infrastructure for large-scale data processing,

increased memory, and bandwidth, as well as a cultural acceptance of

technology in the fabric of our lives. This wasn't true a decade ago."

Cathy O'Neil and Rachel Schutt

Global Average Temperature vs. Number of Pirates

"Correlation is not causation."

"Empirically observed covariation is a necessary but not sufficient condition for causality."

edvard Tufte

"Correlation is not causation but it sure is a hint."
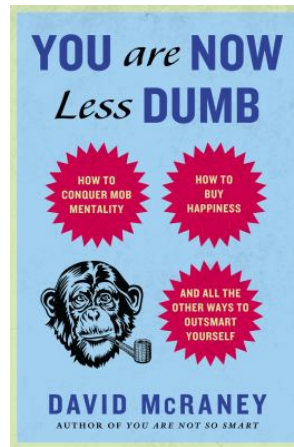
edvard Tufte

## So What?

- Unrelated (Pirates and Climate Change)
- Reverse Causation (Windmills Cause Wind)
- Bi-Directional Causation (Temperature/Pressure)
- Common Causal Variable (Sleeping with your shoes on causes headaches)

"Long before worrying about how to convince others, you first have to understand what's happening yourself."

Andrew Gelman

"Naïve realism, also known as direct realism or common sense realism, is a

philosophy of mind rooted in a theory of perception that claims that the senses

provide us with direct awareness of the external world."

Wikipedia
*http://en.wikipedia.org/wiki/Naïve_realism*

# 1951 Princeton/Dartmouth Game

· Storied rivalry
· Princeton's star player had his nose broken
· Princeton player snapped a Dartmouth player's leg
· Princeton won 13-0
· Editorials from both schools blamed the other
· Two versions of Truth

# They Saw a Game

- Albert Hastorf (Dartmouth) and Hadley Cantril (Princeton) showed the game again to students from both schools
- Asked them to notice infractions, penalties, fill out a questionnaire
- Princeton students 'saw' twice as many infractions by Dartmouth players than Dartmouth students did
- Dartmouth students saw a 'rough but fair' game

"In brief, the data here indicate that there is no such 'thing' as a 'game' existing

'out there' in its own right which people merely 'observe.' The game 'exists' for a

person and is experienced by him only insofar as certain happenings have

significances in terms of his purpose."

Hastorf and Cantril
*They Saw a Game*

"Everything that has ever happened to you has happened inside your skull."

David McRaney
*You Are Now Less Dumb*

## Comparing the Students

- All male
- Ethnic and socioeconomically similar
- Same part of the country
- Same age
- Same basic culture and religious beliefs
- Different schools

"It's a real problem, though, when politicians, CEOs, and other people with the

power to change the way the world works start bungling their arguments for or

against things based on self-delusion generated by imperfect minds and senses."

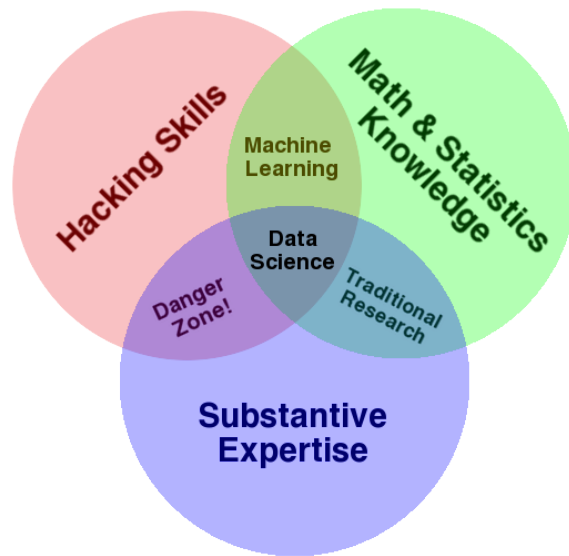David McRaney
*You Are Now Less Dumb*

# Real World Issues

· Foods and Cancer
· Vaccination and Autism
· Global Warming
· GMOs

"Data scientist: n. person who is better at statistics than any software engineer and better at software engineering than any statistician."

Josh Wills

"There's a distinct lack of respect for researchers in academia and industry labs who have been working on this kind of stuff for years, and whose work is based on decades (in some cases, centuries) of work by statisticians, computer scientists, mathematicians, engineers and scientists of all types."

Cathy O'Neil and Rachel Schutt

"Data science is the civil engineering of data. Its acolytes possess a practical knowledge of tools and materials, coupled with a theoretical understanding of what is possible."

Cathy O'Neil and Rachel Schutt

# Data Science Strategy

· Engineering and Infrastructure for collection and logging
· Privacy Access policies
· Role in Decision Making Process

"Narratives are meaning transmitters. They are history-preservation devices. They

create and maintain cultures, and they forge identities that emerge out of the

malleable, imperfect memories of life events."

David McRaney
*You Are Now Less Dumb*

"Your narrative bias makes it nearly impossible for you to really absorb the information from the outside world without arranging it into causes and effects."

David McRaney
*You Are Now Less Dumb*

"Your ancestors invented the scientific method because the common belief fallacy renders your default strategies for making sense of the world generally awful and prone to error."
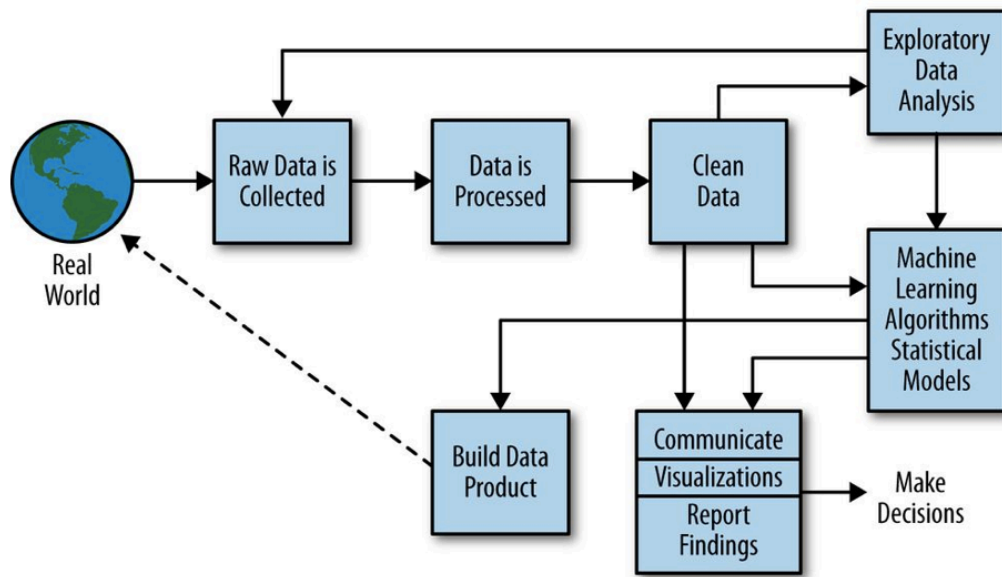
David McRaney
*You Are Now Less Dumb*

菩

# Data Science Techniques

Doing Data Science (O'Neil and Schutt)

# Math

- Statistics
- Linear Algebra
- Numerical Analysis
- Calculus

# Statistics

- Observations and Samples
- Bias
- Modeling
- Distributions
- Fitting a model
- Overfitting

# Techniques

- Linear Regression
- k Nearest Neighbor
- k Means clustering
- Decision Trees

菩

# Programming

## Programming Languages

- C/C++
- Fortran
- Python
- Julia
- R

# Python

- General Purpose, High-Level Programming Language
- Emphasis on readability
- Expressive syntax
- Supports OO, FP, Procedural Programming
- Dynamic
- CPython/Jython

# SciPy

- Ecosystem of open-source packages for science, math and engineering
- NumPy
- SciPy
- Matplotlib
- IPython
- Sympy
- pandas

# SciPy/NumPy

- Numerical analysis
- Optimization problems
- N-Dimensional Arrays
- Linear Algebra
- Fourier transformations

# Julia

- An attempt to create a high-performance, general purpose numerical language
- Think: MATLAB meets Fortran, Python and Lisp
- LLVM-based JIT Compiler
- Growing base of packages
- Impressive performance benchmarks
- MIT License

| | Fortran | Julia | Python | R | Matlab | Octave | Mathe-matica | JavaScript | Go |
|---|---|---|---|---|---|---|---|---|---|
| | gcc 4.8.1 | 0.2 | 2.7.3 | 3.0.2 | R2012a | 3.6.4 | 8.0 | V8 3.7.12.22 | go1 |
| fib | 0.26 | 0.91 | 30.37 | 411.36 | 1992.00 | 3211.81 | 64.46 | 2.18 | 1.03 |
| parse_int | 5.03 | 1.60 | 13.95 | 59.40 | 1463.16 | 7109.85 | 29.54 | 2.43 | 4.79 |
| quicksort | 1.11 | 1.14 | 31.98 | 524.29 | 101.84 | 1132.04 | 35.74 | 3.51 | 1.25 |
| mandel | 0.86 | 0.85 | 14.19 | 106.97 | 64.58 | 316.95 | 6.07 | 3.49 | 2.36 |
| pi_sum | 0.80 | 1.00 | 16.33 | 15.42 | 1.29 | 237.41 | 1.32 | 0.84 | 1.41 |
| rand_mat_stat | 0.64 | 1.66 | 13.52 | 10.84 | 6.61 | 14.98 | 4.52 | 3.28 | 8.12 |
| rand_mat_mul | 0.96 | 1.01 | 3.41 | 3.98 | 1.10 | 3.41 | 1.16 | 14.60 | 8.51 |

**Figure:** benchmark times relative to C (smaller is better, C performance = 1.0).
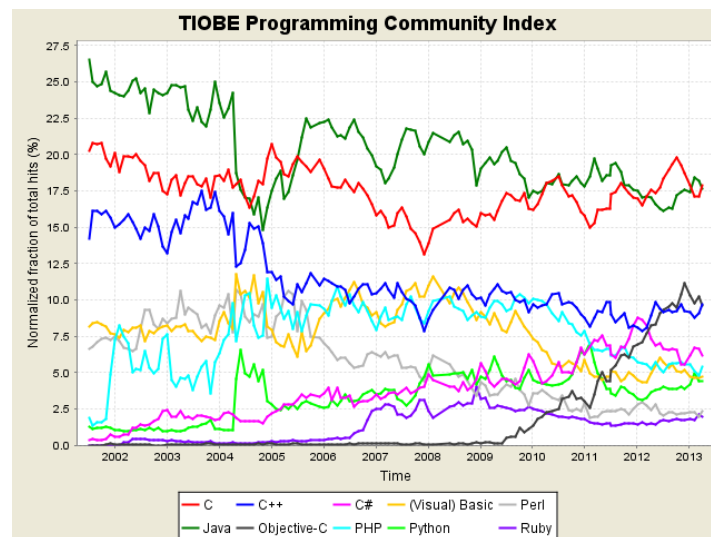
# Nerdy Language Features

· Multiple dispatch
· Dynamic type system
· Built-in package manager
· Lisp-like macros and other metaprogramming fu
· Ability to call C/Python functions
· Supports parallel and distributed processing
· Coroutines

# R

- Created by Ross Ihaka and Robert Gentleman
- Maintained by the R Development Core Team
- Part of the GNU Project
- Commercial support
- Implemented in C, Fortran and R

TIOBE Programming Community Index

| Position Apr 2013 | Position Apr 2012 | Delta in Position | Programming Language | Ratings Apr 2013 | Delta Apr 2012 | Status |
|---|---|---|---|---|---|---|
| 1 | 1 | = | C | 17.862% | +0.31% | A |
| 2 | 2 | = | Java | 17.681% | +0.65% | A |
| 3 | 3 | = | C++ | 9.714% | +0.82% | A |
| 4 | 4 | = | Objective-C | 9.598% | +1.36% | A |
| 5 | 5 | = | C# | 6.150% | -1.20% | A |
| 6 | 6 | = | PHP | 5.428% | +0.14% | A |
| 7 | 7 | = | (Visual) Basic | 4.699% | -0.26% | A |
| 8 | 8 | = | Python | 4.442% | +0.78% | A |
| 9 | 10 | ↑ | Perl | 2.335% | -0.05% | A |
| 10 | 11 | ↑ | Ruby | 1.972% | +0.46% | A |
| 11 | 9 | ↓↓ | JavaScript | 1.509% | -1.37% | A |
| 12 | 14 | ↑↑ | Visual Basic .NET | 1.095% | +0.12% | A |
| 13 | 15 | ↑↑ | Lisp | 0.905% | -0.05% | A |
| 14 | 16 | ↑↑ | Pascal | 0.887% | +0.07% | A |
| 15 | 13 | ↓↓ | Delphi/Object Pascal | 0.840% | -0.53% | A |
| 16 | 32 | ↑↑↑↑↑↑↑↑↑↑ | Bash | 0.840% | +0.47% | A |
| 17 | 18 | ↑ | Transact-SQL | 0.723% | -0.04% | A |
| 18 | 12 | ↓↓↓↓↓↓ | PL/SQL | 0.715% | -0.66% | A |
| 19 | 24 | ↑↑↑↑↑ | Assembly | 0.710% | +0.24% | A-- |
| 20 | 21 | ↑ | Lua | 0.650% | +0.08% | B |

| Position | Programming Language | Ratings |
|---|---|---|
| 21 | Ada | 0.642% |
| 22 | SAS | 0.634% |
| 23 | ABAP | 0.588% |
| 24 | MATLAB | 0.517% |
| 25 | COBOL | 0.491% |
| 26 | R | 0.484% |
| 27 | Scheme | 0.419% |
| 28 | Fortran | 0.407% |
| 29 | Scala | 0.336% |
| 30 | Prolog | 0.324% |
| 31 | Erlang | 0.323% |
| 32 | Haskell | 0.317% |
| 33 | Scratch | 0.317% |
| 34 | Logo | 0.316% |
| 35 | D | 0.314% |
| 36 | Forth | 0.240% |
| 37 | Smalltalk | 0.235% |
| 38 | ActionScript | 0.226% |
| 39 | APL | 0.222% |
| 40 | Common Lisp | 0.219% |

## Supports
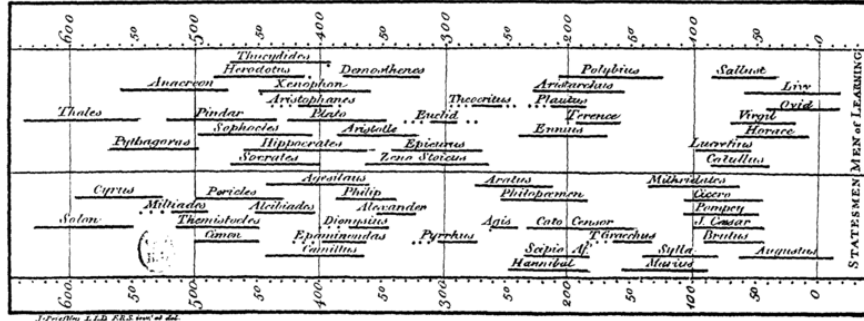
· Basic Math
· Statistical Analysis
· Optimization Problems
· Signal Processing
· Graphics and Visualization
· Data Mining
· Machine Learning

菩

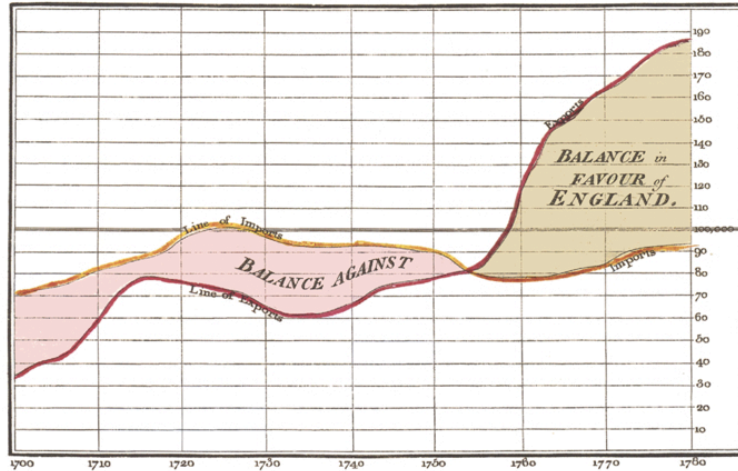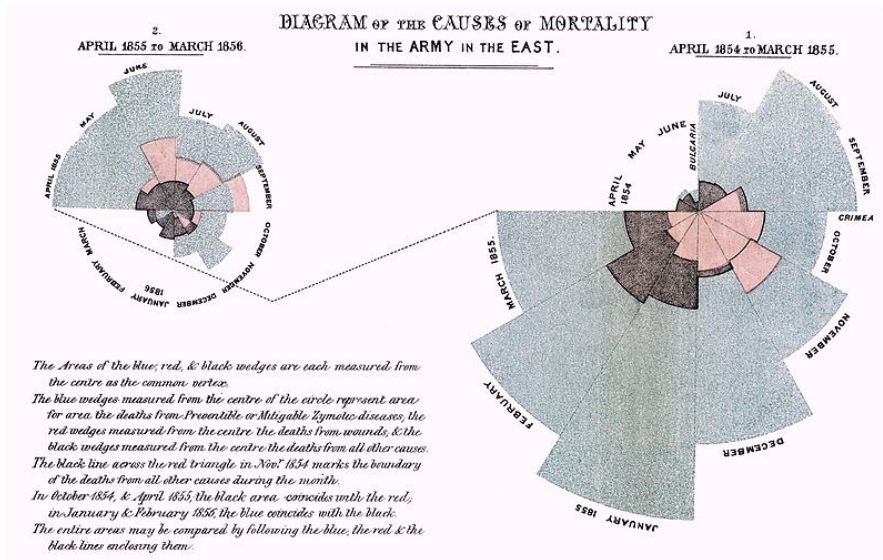# Visualization

A Specimen of a Chart of Biography.

Exports and Imports to and from DENMARK & NORWAY from 1700 to 1780.

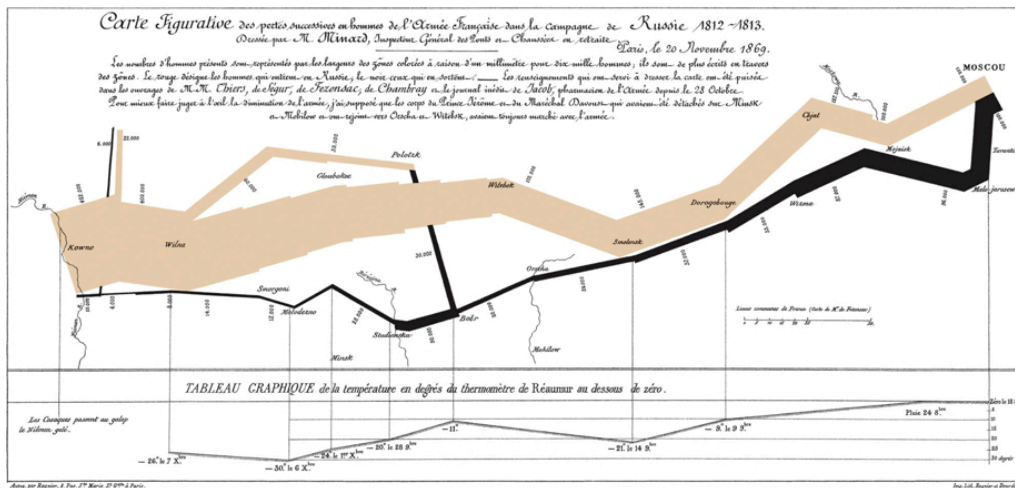The Bottom line is divided into Years, the Right hand line into L10,000 each.

CHART Representing the EXTENT, POPULATION & REVENUES, of the PRINCIPAL NATIONS in EUROPE, after the DIVISION of POLAND & TREATY of LUNEVILLE.

Yards
50   0   50   100   150   200

✕ Pump      • Deaths from cholera
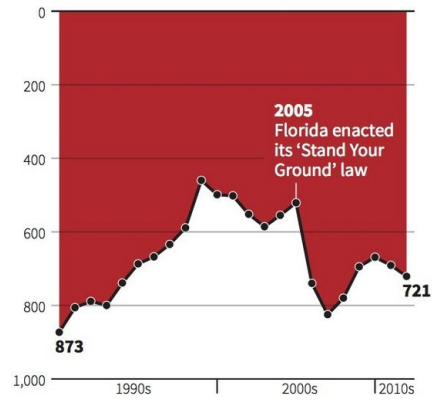
History of O-Ring Damage in Field Joints

# Techniques

- Graphical Analysis
- Presentation Graphics

# Gun deaths in Florida

Number of murders committed using firearms



**2005**
Florida enacted its 'Stand Your Ground' law

**873**

**721**

Source: Florida Department of Law Enforcement

C. Chan 16/02/2014

REUTERS
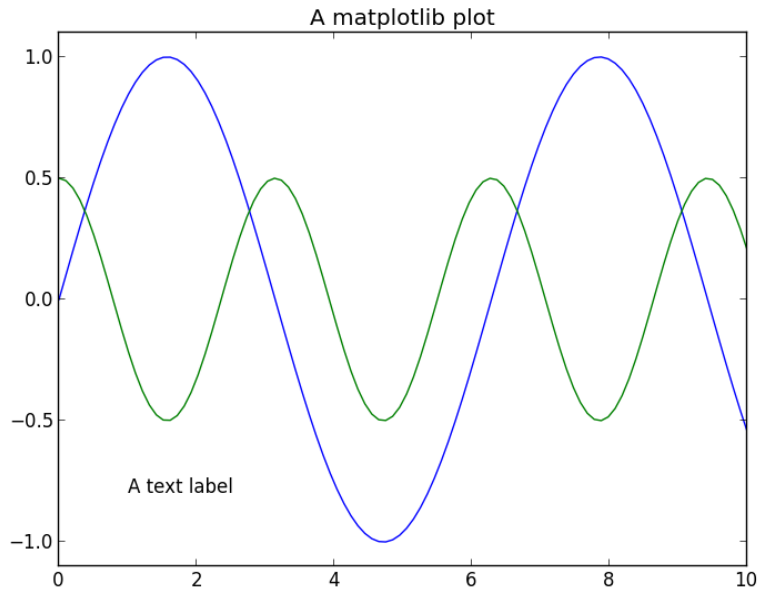
```
ipython -pylab

x = linspace( 0, 10, 100 )
plot( x, sin(x) )
plot( x, 0.5*cos(2*x) )
title( "A matplotlib plot" )
text( 1, -0.8, "A text label" )
ylim( -1.1, 1.1 )
savefig( 'matplotlib.png' )
```

A matplotlib plot

# d3.js

- Data-Driven Documents
- JavaScript library that uses HTML, SVG and CSS
- Bind data to the DOM
- Data-driven transformations

https://github.com/mbostock/d3/wiki/Gallery

http://mbostock.github.io/d3/talk/20111116/iris-splom.html

http://bost.ocks.org/mike/uberdata/

http://exposedata.com/parallel/

http://mbostock.github.io/d3/tutorial/circle.html

菩

# Machine Learning

菩

Data Mining

菩

Big Data

"First, it is a bundle of technologies. Second, it is a potential revolution in measurement. And third, it is a point of view, or philosophy, about how decisions will be— and perhaps should be— made in the future."

Steve Lohr, New York Times (2013-10-09)

## Everything You Know About Something

| ID | Col1 | Col2 | Col3 | Col4 | Col5 | Col6 | .... | ColN |
|---|---|---|---|---|---|---|---|---|
| Thing1 | Value1 | Value2 | Value3 | Value4 | Value5 | Value6 | .... | ValueN |

# Everything You Know About Everything

| ID | Col1 | Col2 | Col3 | Col4 | Col5 | Col6 | .... | ColN |
|---|---|---|---|---|---|---|---|---|
| **Thing1** | Value1 | Value2 | Value3 | | Value5 | | .... | ValueN |
| **Thing2** | Value1 | | Value3 | Value4 | Value5 | Value6 | .... | ValueN |
| **Thing3** | | Value2 | Value3 | | Value5 | Value6 | .... | ValueN |
| **Thing4** | Value1 | Value2 | Value3 | Value4 | Value5 | Value6 | .... | ValueN |
| ... | ... | ... | ... | ... | ... | ... | .... | ... |

# Distribute Rows in their Entirety

| ID | Col1 | Col2 | Col3 | Col4 | Col5 | Col6 | .... | ColN |
|---|---|---|---|---|---|---|---|---|
| **Thing1** | Value1 | Value2 | Value3 | | Value5 | | .... | ValueN |
| **Thing3** | | Value2 | Value3 | | Value5 | Value6 | .... | ValueN |

# Distribute Columns in their Entirety

| ID | Col2 | Col3 | Col5 | ColN |
|---|---|---|---|---|
| Thing1 | Value2 | Value3 | | ValueN |
| Thing3 | Value2 | Value3 | Value5 | ValueN |
| Thing4 | Value2 | Value3 | Value5 | ValueN |
| ... | ... | ... | ... | ... |

菩

# Linked Data

# Distribute Arbitrary Cells

| ID | Col2 | Col3 | Col5 | ColN |
|---|---|---|---|---|
| Thing1 | | Value3 | | ValueN |
| Thing3 | | | Value5 | ValueN |
| Thing4 | Value2 | Value3 | | ValueN |
| … | … | … | … | … |

# Linking Open Data Project

· Started in 2007 by W3C Semantic Web Education and Outreach(SWEO) Interest Group
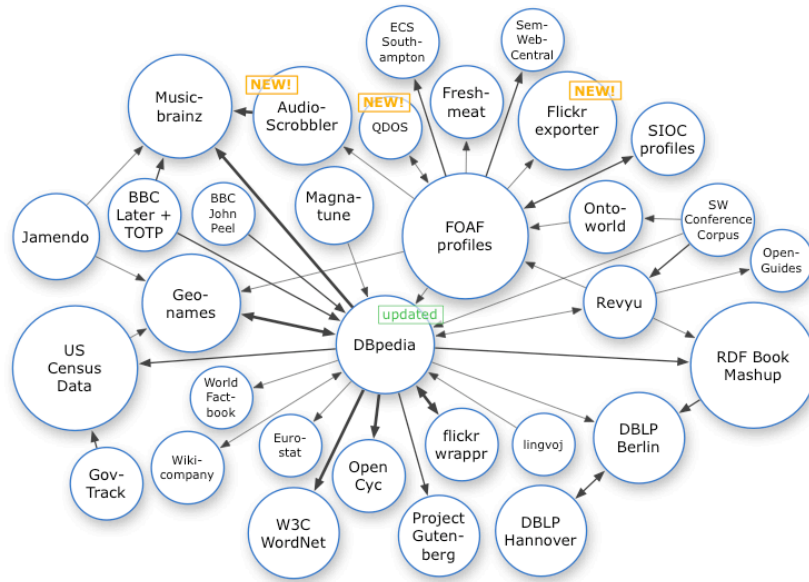· Make data freely available
· Doubled in size every 10 months

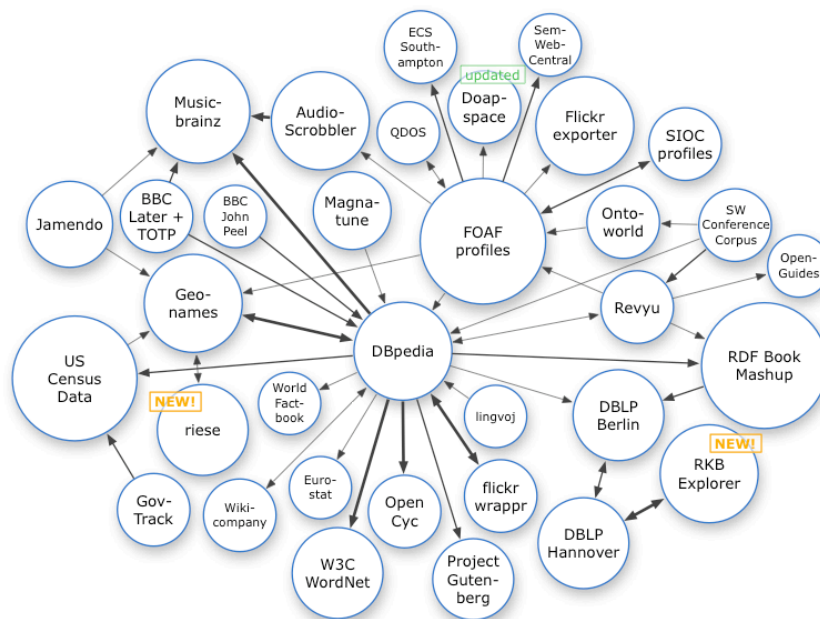As of May 2007
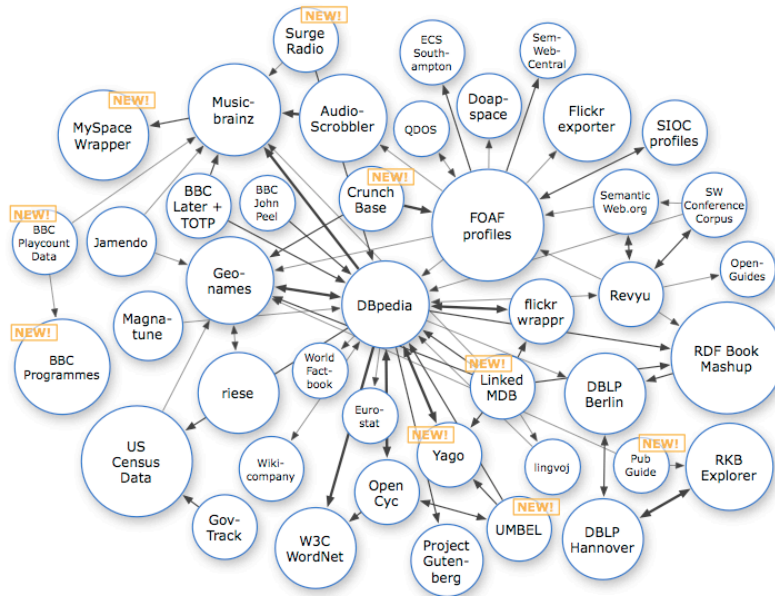
As of September 2008

As of March 2009

As of March 2009

Media
Geographic
Publications
User-generated content
Government
Cross-domain
Life sciences

As of September 2010

| Domain | # Datasets | # Triples | # Links |
|---|---|---|---|
| Media | 25 | 1,800,000,000 | 50,000,00 |
| Geographic | 31 | 6,000,000,000 | 35,000,000 |
| Government | 49 | 13,000,000,000 | 19,000,000 |
| Publications | 87 | 2,900,000,000 | 140,000,000 |
| Cross-Domain | 41 | 4,100,000,000 | 63,000,000 |
| Life Sciences | 41 | 3,000,000,000 | 191,000,000 |
| User-Generated Content | 20 | 134,000,000 | 3,400,000 |
| Total | 295 | 31,000,000,000 | 504,000,000 |

http://lod-cloud.net/state/

## Domains By Datasets

- Media
- Geographic
- Government
- Publications
- Cross-Domain
- Life Sciences
- User-Generated Content

Media 9%
Geographic 11%
Government 17%
Publications 30%
Cross-Domain 14%
Life Sciences 14%
User-Generated Content 7%

## Domains By Triples

- Media
- Life Sciences
- Cross-Domain
- User-Generated Content
- Publications
- Government
- Geographic

Media 6%
Life Sciences 10%
Cross-Domain 13%
User-Generated Content 0%
Publications 9%
Government 42%
Geographic 19%

Domains By Links



Media
Cross-Domain
Publications
Government
Geographic
Life Sciences
User-Generated Content

10%
13%
28%
4%
7%
38%
1%

# DBPedia

- Linked Dataset derived from Wikipedia
- Creative Commons Attribution-ShareAlike 3.0 License
- GNU Free Documentation License
- Multi-domain
- Consensus-based
- Kept current by Wikipedia activity
- Multi-lingual

# DBPedia Numbers (English Version)

http://dbpedia.org/About

- Describes 4 million things
- 3.22 million are classified by an ontology
- 832,000 people
- 639,000 places
- 372,000 creative works
- 209,000 organizations
- 226,000 species
- 5,600 diseases

# DBPedia Numbers (Non-English Version)

http://dbpedia.org/About

- 119 Localized Language Versions
- Describe 24.9 million things (w/ repetition)
- 16.8 million are connected to English DBPedia

# DBPedia Summary

http://wiki.dbpedia.org/Datasets39/DatasetStatistics?v=dqp

- Overall 12.6 million unique things
- 24.6 million links to images
- 27.6 million links to pages
- 45 million links to other RDF datasets
- 67 million links to Wikipedia categories
- 41.2 million links to YAGO categories
- 2.46 billion RDF triples
- 470 million (English), 1.98 billion (Non-English)

# Use Cases

http://wiki.dbpedia.org/UseCases?v=ene

- Improve Wikipedia Search
- Include DBPedia data in your documents
- Support for Geographic Data
- Documentation Classification, Annotation
- Multi-Domain Ontology

DBPedia

http://dbpedia.org

Most Important Query Ever Run

http://tinyurl.com/n9hhs68

http://www.r-bloggers.com/sparql-with-r-in-less-than-5-minutes/

http://linkedscience.org/tools/sparql-package-for-r/tutorial-on-sparql-package-for-r/

菩

# Books

# Data Science Books

· Data Analysis w/ Open Source Tools, Philipp K. Janet (ORA)
· Data Smart: Using Data Science to Transform Information into Insight, John W. Foreman (Wiley)
· Doing Data Science : Straight Talk from the Frontline, Cathy O'Neil, Rachel Schutt (ORA)
· The R Book, Michael J. Crawley (Wiley)
· R Tutorial w/ Bayesian Statistics Using OpenBUGS, Chi Yau
· Applied Predictive Modeling, Max Kuhn and Kjell Johnson (Springer)
· Introductory Statistics w/ R, Peter Dalgaard (Wiley)
· Think Stats, Allen B. Downey (ORA)
· R Cookbook, Paul Teetor (ORA)
· R Graphics Cookbook, Winston Chang (ORA)

HTML

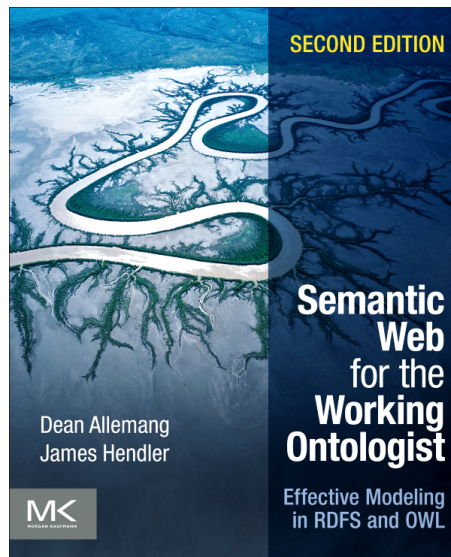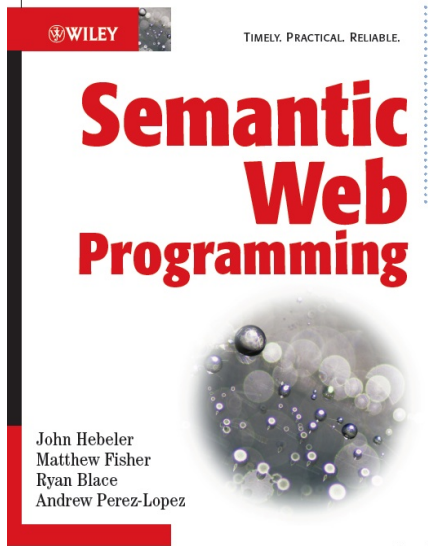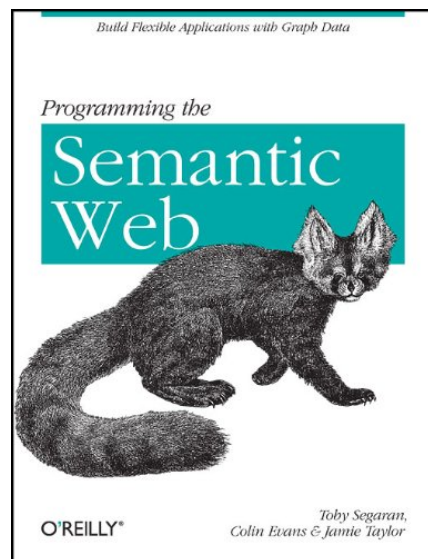http://www.quora.com/Data-Science/How-do-I-become-a-data-scientist

HTML

http://cm.bell-labs.com/cm/ms/departments/sia/doc/datascience.pdf

**Semantic Web Programming**

John Hebeler
Matthew Fisher
Ryan Blace
Andrew Perez-Lopez

*Build Flexible Applications with Graph Data*

*Programming the*

**Semantic Web**

O'REILLY®

*Toby Segaran,
Colin Evans & Jamie Taylor*

MORGAN&CLAYPOOL PUBLISHERS

# Linked Data
*Evolving the Web into a
Global Data Space*

**Tom Heath**

**Christian Bizer**

# Linked Data
Structured data on the Web

David Wood
Marsha Zaidman
Luke Ruth
with Michael Hausenblas

FOREWORD BY Tim Berners-Lee

MANNING

Resource-Oriented Architecture Patterns for Webs of Data

Brian Sletten

## Questions?

✉ brian@bosatsu.net

🐦 @bsletten

g+ http://tinyurl.com/bjs-gplus

○ bsletten