

A drink from the  
fire hose so you can  
start your own  
exploration.

# Introduction to Data Science

Bryan Nehl  
[@k0emt](https://twitter.com/k0emt)

[dbBear.com](http://dbBear.com)

# What is Data Science?

0Turns out, that's a  
good question!

0Skills

0Sourcing

0Extract, Transform  
and Load

0Databases

0Languages and  
Libraries

0Business Knowledge

0Visualization

0What now?

# Skills

## 0 Data Journal – Engineer's Notebook

0 <http://soloso.blogspot.com/2014/05/engineers-notebook.html>

## 0 Regular Expressions, Markdown

## 0 Analysis (math/stats is part of this!)

## 0 LINUX

0 Putty

## 0 Version Control

0 <http://git-scm.com/>

## 0 Messaging – RabbitMQ

0 <http://rabbitmq.com>



# Virtualization

0 Why?

0 Time

0 Pre-built images

0 Cost

0 On Demand

0 How/Where?

0 Microsoft Azure

0 Amazon Elastic Cloud

0 Google Compute

0 Heroku

0 Your Own Machine

0 Oracle VirtualBox -

[www.virtualbox.org](http://www.virtualbox.org)

0 Parallels

0 VMWare

# Common Virtual Machines (VMs)

LAMP/WAMP

0 Linux/Windows

0 Apache

0 MySQL

0 PHP

MEAN

0 MongoDB

0 Express

0 Angular

0 Node.JS

0 mean.io



MongoDB is the leading NoSQL database, empowering businesses to be more agile and scalable.

express

Express is a minimal and flexible node.js web application framework, providing a robust set of features for building single and multi-page, and hybrid web applications.



AngularJS lets you extend HTML vocabulary for your application. The resulting environment is extraordinarily expressive, readable, and quick to develop.



Node.js is a platform built on Chrome's JavaScript runtime for easily building fast, scalable network applications.

# Data Sourcing

# Sourcing the data

0Locate it

0Provided

0Search for it

0Manually

0Automated

0Networking

0Get it

0ftp

0Scraping

0Database

0Web services

0Work with it

# Some dataz

0 data.gov

0 data.mo.gov

0 data.kcmo.gov

0 data.gov.uk

0 data.worldbank.org

0 gutenberg.org

0 [http://bitly.com/  
bundles/hmason/1](http://bitly.com/bundles/hmason/1)

0 [soloso.blogspot.com  
/2011/07/getting-  
enron-mail-  
database-into.html](http://soloso.blogspot.com/2011/07/getting-enron-mail-database-into.html)

# ETL

Extract, Transform and Load (the data!)

# Extract, Transform & Load

## Data Formats

0 XLS

0 CSV

0 Text

0 Delimited

0 Fixed format

0 JSON – json.org

0 XML & HTML

0 Mail files

0 SQL scripted INSERTs

0 PDF

0 Character sets

<https://docs.python.org/2/howto/unicode.html>

# Data Analysis ???

- 0 What do I expect to see?
- 0 What are the field types?
- 0 Does the field type change?
- 0 What are the range of values?
- 0 How frequently do those values occur?
  - 0 *Can I get a graph please?*
- 0 Are there nulls?
- 0 How big is my sample set?
  - 0 *Is it significant?*
- 0 How big do I expect the real data to be?
- 0 Are there holes in the data?
- 0 What constitutes a *good* record?
- 0 Where are the trends/ clusters in the data?



0 [openrefine.org](http://openrefine.org)

0 Explore

0 Clean and Refine

0 Matching

0 Exporting

0 Plug-in architecture

0 Works with web services

0 Works with [freebase.com](http://freebase.com)

# Languages and Libraries

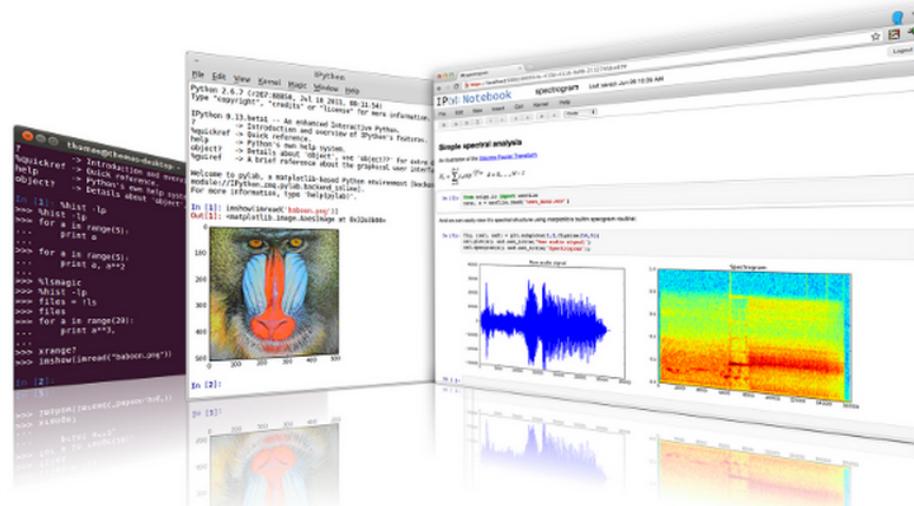
Tools to get the work done. Don't reinvent the wheel.

# Languages

0 Python – [python.org](http://python.org)  
0 IPython – [ipython.org](http://ipython.org)



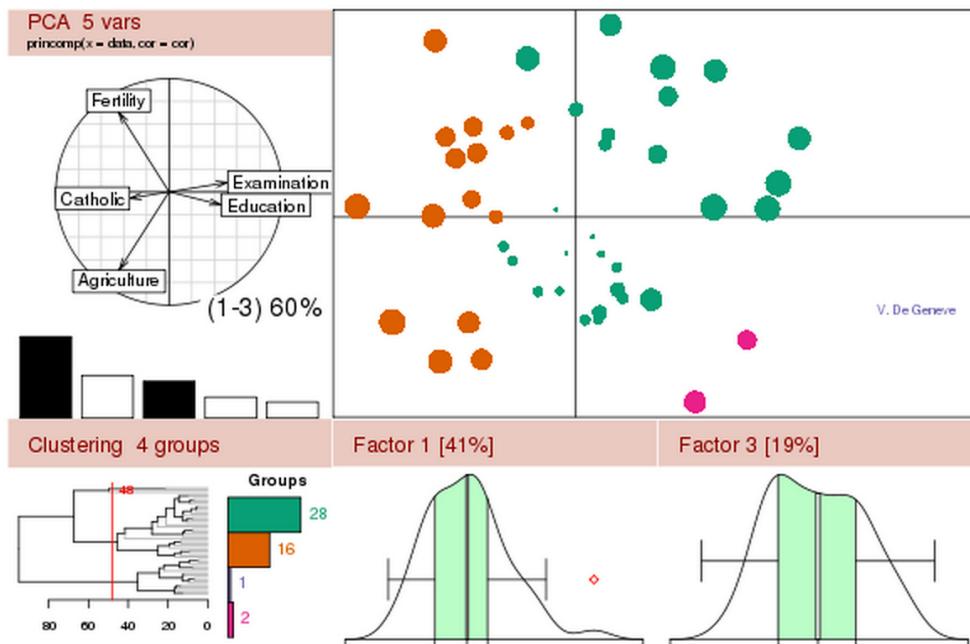
IP[y]: IPython  
Interactive Computing



# Languages

0 R - [www.r-project.org](http://www.r-project.org)

0 R Studio - [www.rstudio.com](http://www.rstudio.com)



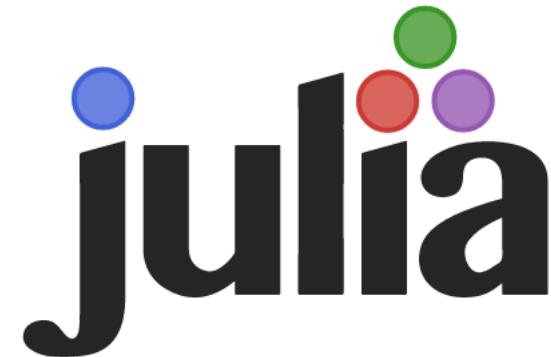
# Languages

0 Julia – [julialang.org](http://julialang.org)

- 0 Very fast numeric operations
- 0 Close to FORTRAN in speed
- 0 Designed for parallelism
- 0 IJulia works with IPython

0 Javascript

- 0 jquery



# Libraries for Excel

0 It is everywhere

0 Python Libraries:

0 xlrd

0 XlsxWriter

0 Apache Project – Office Open XML file formats

0 <http://poi.apache.org/>

0 Excel

0 Word

0 PowerPoint

# Libraries

0 SciPy

0 [scipy.org](http://scipy.org)

0 NumPy

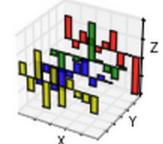
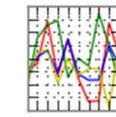
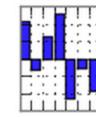
0 [numpy.org](http://numpy.org)

0 Pandas - Python data analysis library

0 [pandas.pydata.org](http://pandas.pydata.org)



pandas  
 $y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$



# Libraries

0 lxml

0 [lxml.de](http://lxml.de)

0 pymongo

0 [pypi.python.org/pypi/pymongo/](http://pypi.python.org/pypi/pymongo/)

0 pika – AMQP

0 [pypi.python.org/pypi/pika](http://pypi.python.org/pypi/pika)

0 nose – unit test framework extension

0 [nose.readthedocs.org](http://nose.readthedocs.org)



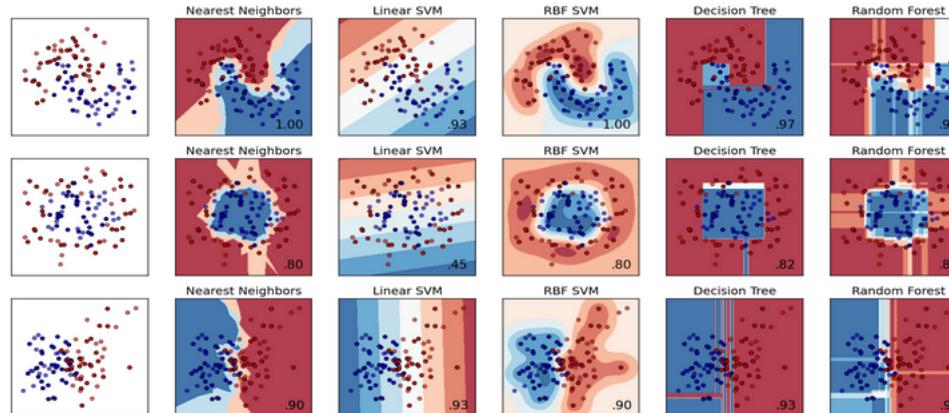
# scikit-learn.org

0 Machine Learning

0 Clustering

0 Classification

0 Data Mining



# Packages

0 **Anaconda** Scientific Python development environment

0 Getting IPython set up by hand is a pain—Anaconda is a must on Windows machines.

0 <https://store.continuum.io/cshop/anaconda/>

0 **wakari.io** web based Python data analysis

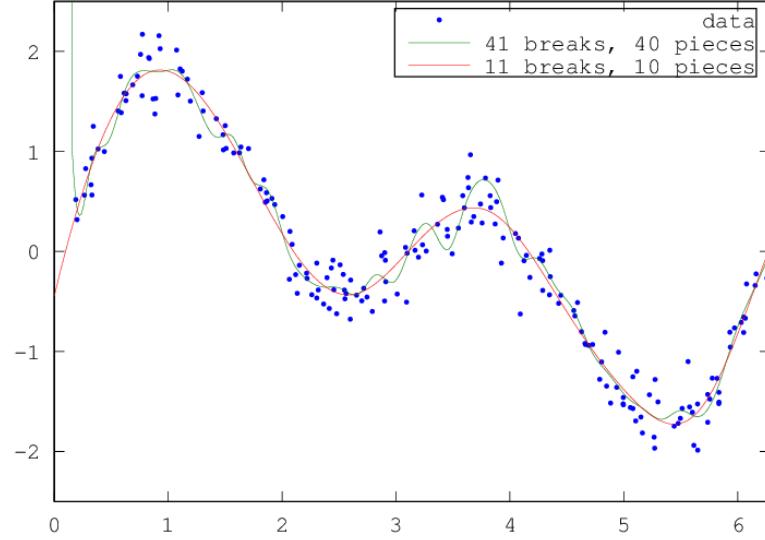


# Other Tools

0 Processing.org

0 Octave

0 [www.gnu.org/software/octave/](http://www.gnu.org/software/octave/)



# Databases

Choose the right one(s) for the job!

Polyglot Persistence

<http://martinfowler.com/bliki/PolyglotPersistence.html>

# Relational - SQL

0 MySQL

0 open source

0 Oracle

0 Microsoft SQL Server

0 Express Editions

0 Microsoft Access

0 ODBC / JDBC



# NOSQL

## 0 Definition

### 0 MongoDB – JSON/BSON documents

0 [mongodb.org](http://mongodb.org)

### 0 neo4j – graph

0 [neo4j.org](http://neo4j.org)

### 0 PostgreSQL – object-relational

0 [postgresql.org](http://postgresql.org)



# NOSQL

0 Aerospike – Key-Value

0 [aerospike.com](http://aerospike.com)



IN-MEMORY NoSQL DATABASE

0 redis – Key-Value

0 [redis.io](http://redis.io)



0 CouchDB – JSON docs, HTTP

0 [couchdb.apache.org](http://couchdb.apache.org)



0 Couchbase – JSON, memcached

0 [couchbase.com](http://couchbase.com)

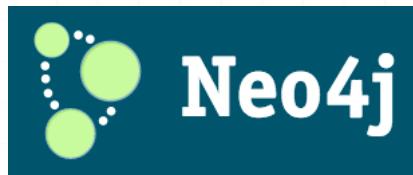




- 0 Distributed framework for processing large datasets
  - 0 MongoDB and other databases can be used to feed it
- 
- 0 [hadoop.apache.org](http://hadoop.apache.org)

# Example Polyglot System

## Lobbyist Influence

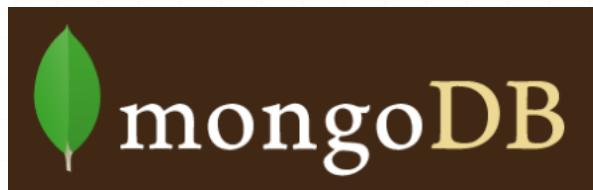


0 Relationships:

0 Politicians

0 Lobbyists

0 Legislation



0 PDF Versions of the bill



0 Financial information

# Business Knowledge

I have data. Now what?

# Numbers need context

Visitors

1M

Page Views

5.2M

72%

Conversion Rate

42

Customer average age

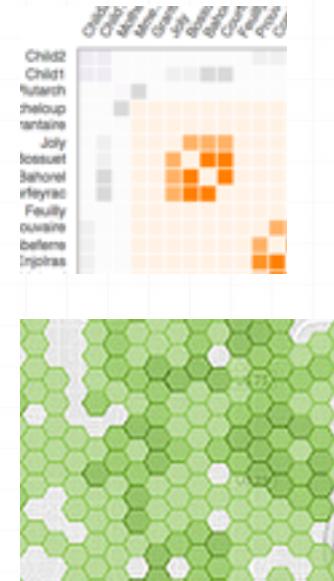
1

Top Referrer.com

# Analysis

## Techniques

- 0 Adjacency Matrix
- 0 pivot and fold operations on tables
- 0 hexagonal binning
- 0 confusion matrix
- 0 predictive modeling fundamentals
- 0 machine learning
- 0 The work of **John Tukey** (Statistics)
  - 0 [wikipedia.org/wiki/John\\_Tukey](https://en.wikipedia.org/wiki/John_Tukey)



# What can I do with this data that will benefit the business?

- 0 Is there some insight I can bring?
- 0 Can I generalize from this data? (global)
- 0 Can I ascertain local area insights?
- 0 Are there natural partitions in the data?
  - 0 Gender, race, age, location?
- 0 Is there some business pain I can relieve?
- 0 Can I enhance an existing data set?
- 0 Can I bring in the data product with a shorter cycle time?

# Business Intelligence

0 [wikipedia.org/wiki/Business\\_intelligence](https://wikipedia.org/wiki/Business_intelligence)

0 Pentaho - [pentaho.com](http://pentaho.com)

0 Tableau - [tableausoftware.com](http://tableausoftware.com)

0 Microsoft Power BI for Office 365

0 [microsoft.com/en-us/powerbi/](http://microsoft.com/en-us/powerbi/)



# Visualization

Use your pixels!

# HTML Tools & Libraries

0 HTML5 / CSS3

0 Javascript

0 **d3.js – d3js.org**

0 HTML 5 canvas charts

0 chartjs.org

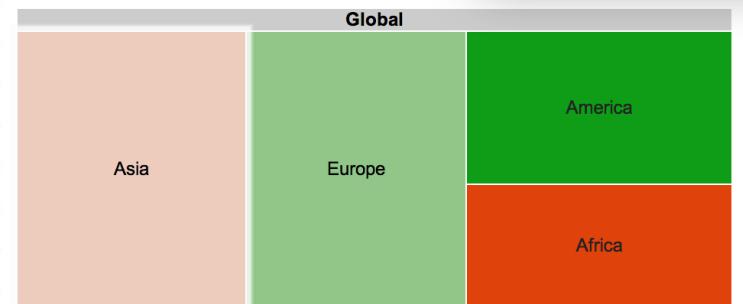
0 canvasjs.com



# Google Tools & Libraries

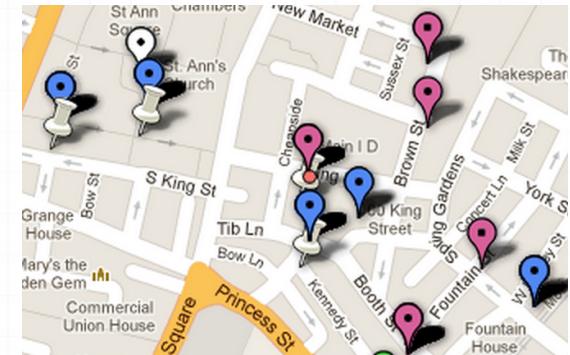
0 Google Charts

0 [developers.google.com/chart/](http://developers.google.com/chart/)



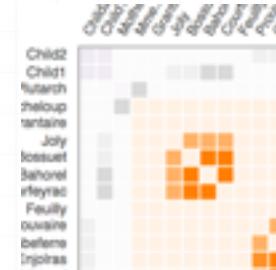
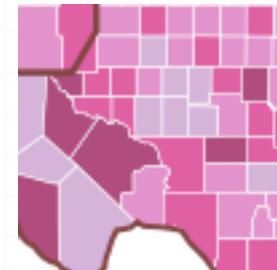
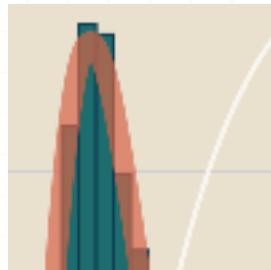
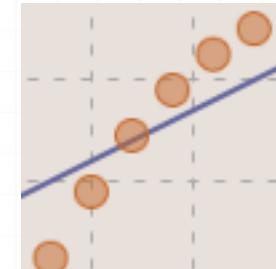
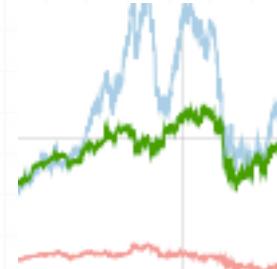
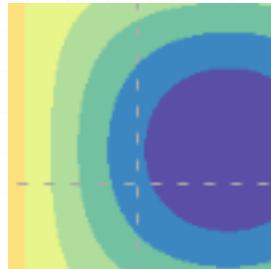
0 Google Fusion Tables

0 Now integrated with Google Drive



# Python Tools & Libraries

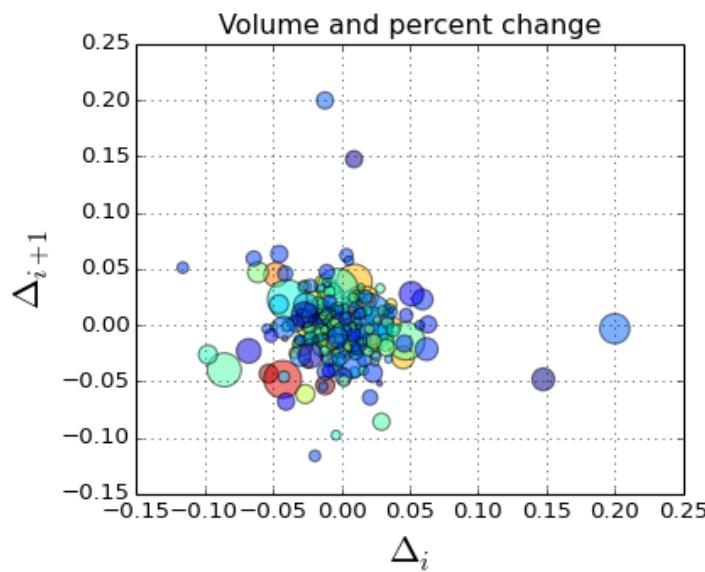
0 bokeh - [bokeh.pydata.org](http://bokeh.pydata.org)



# Python Tools & Libraries

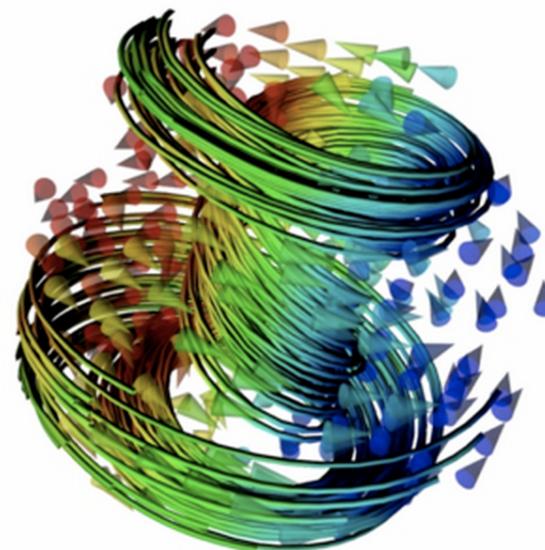
## 0 Matplotlib

0 [matplotlib.org](http://matplotlib.org)



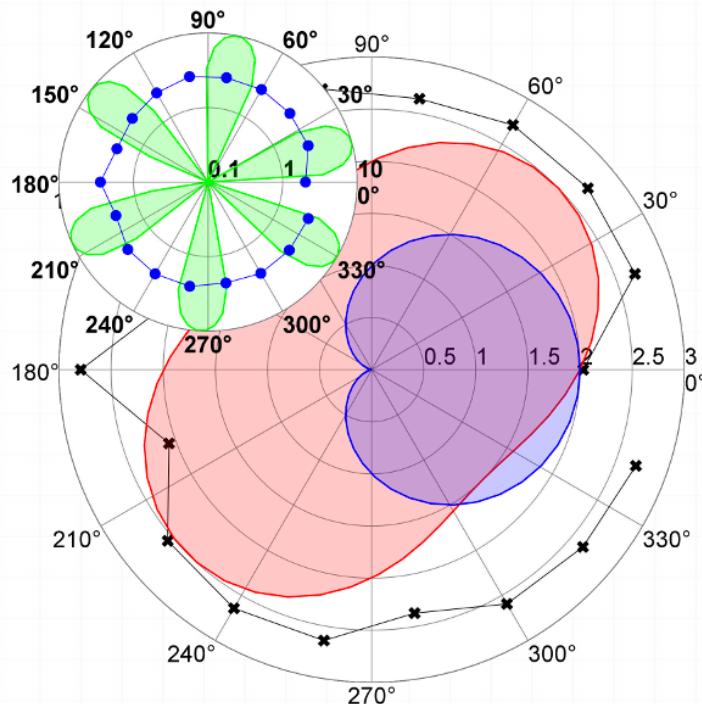
## 0 Mayavi 2

0 [code.enthought.com/projects/mayavi/](http://code.enthought.com/projects/mayavi/)

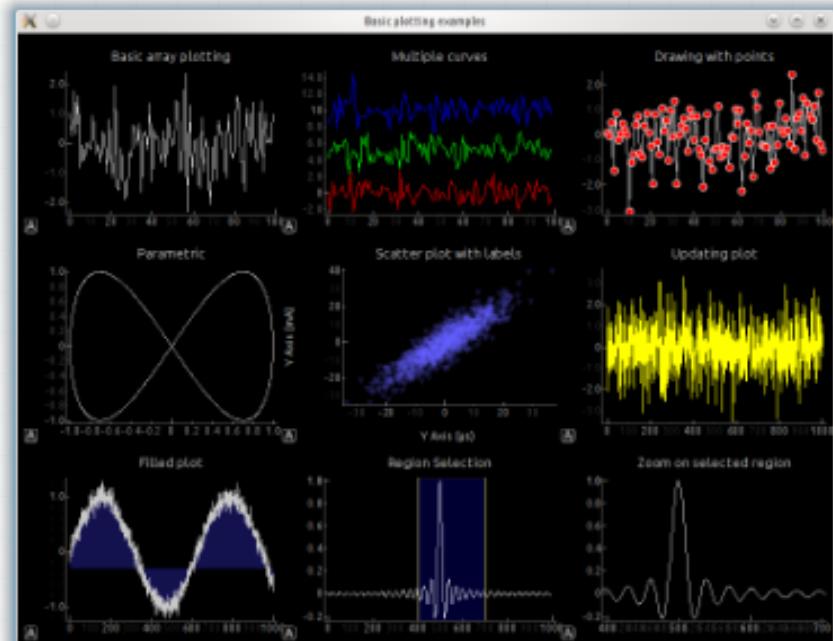


# Python Tools & Libraries

0 [home.gna.org/veusz/](http://home.gna.org/veusz/)



0 [pyqtgraph.org](http://pyqtgraph.org)



# Excel: data quality check

0.335857	0.733451	0.599874	0.335857	0.733451	0.599874
0.398299	0.193938	0.572766	0.398299	0.193938	0.572766
0.71445	0.22316	0.360831	0.71445	0.22316	0.360831
0.821805	0.568467	0.858095	0.821805	0.568467	0.858095
0.069867	0.434296	0.730381	0.069867	0.434296	0.730381
0.206457	0.918653	0.377569	0.206457	0.918653	0.377569
0.04397	0.908735	0.801125	0.04397	0.908735	0.801125
0.952784	0.213182	0.621818	0.952784	0.213182	0.621818
0.305901	0.528717	0.545583	0.305901	0.528717	0.545583
0.732739	0.579152	0.202078	0.732739	0.579152	0.202078

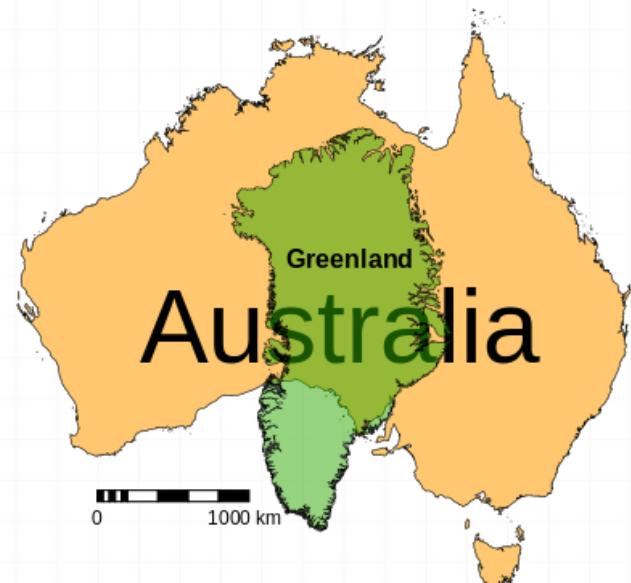
# Map Projections

0 [http://en.wikipedia.org/wiki/List\\_of\\_map\\_projections](http://en.wikipedia.org/wiki/List_of_map_projections)

0 <http://xkcd.com/977/>

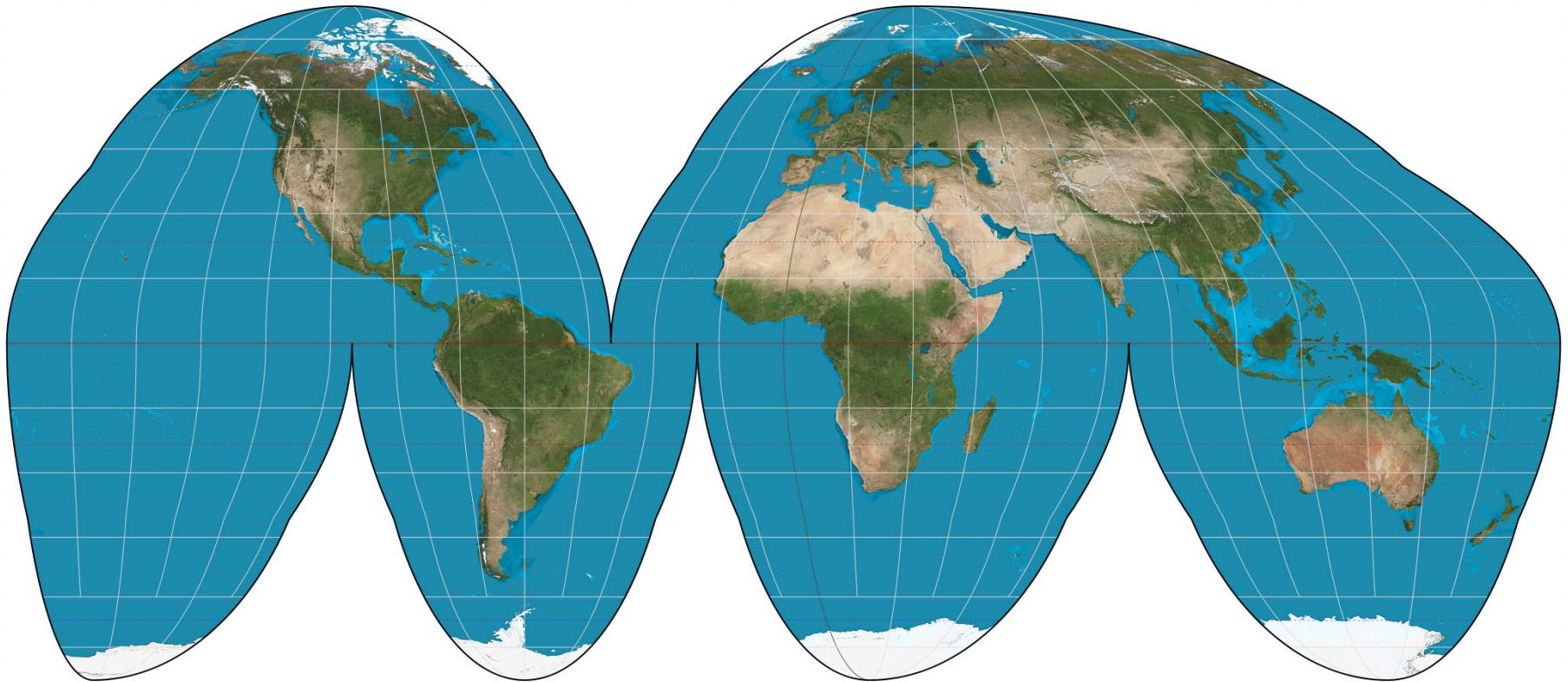
0 <http://www.jasondavies.com/maps/countries-by-area/>

# Mercator



[http://upload.wikimedia.org/wikipedia/commons/f/f4/Mercator\\_projection\\_SW.jpg](http://upload.wikimedia.org/wikipedia/commons/f/f4/Mercator_projection_SW.jpg)  
[http://commons.wikimedia.org/wiki/File:Australia-Greenland\\_size\\_comparison.svg](http://commons.wikimedia.org/wiki/File:Australia-Greenland_size_comparison.svg)

# Goode homolosine projection



Equal area projection map

[http://en.wikipedia.org/wiki/File:Goode\\_homolosine\\_projection\\_SW.jpg](http://en.wikipedia.org/wiki/File:Goode_homolosine_projection_SW.jpg)

*That was a  
**BIG** drink  
from the fire hose!*

What now?



# Teams

# Decide your direction

0 Personal Tech Radar

[http://nealford.com/memeagora/2013/05/28/  
build\\_your\\_own\\_technology\\_radar.html](http://nealford.com/memeagora/2013/05/28/build_your_own_technology_radar.html)

# Training

- 0 Tutorials and sample files that come with software.
- 0 Local courses
- 0 Online Education from vendors
  - 0 MongoDB University  
[0 university.mongodb.com](http://university.mongodb.com)
  - 0 Other online education
    - 0 Coursera – [coursera.org](http://coursera.org)
    - 0 iTunes University (iTunes U)  
[0 oreilly.com/data/free](http://oreilly.com/data/free)
    - 0 O'Reilly Safari, books, videos and free publications  
[0 oreilly.com/data/free](http://oreilly.com/data/free)

# Experiment

- 0 Set up a development environment
- 0 Create a Virtual Machine
- 0 Try out stuff
  - 0 Work related
  - 0 Something you are passionate about
- 0 Share your experiences
  - 0 blog, tweet, present
  - 0 GitHub and Gists

# Contests

0 kaggle.com

0 www.kdnuggets.com/competitions/

0 www.crowdanalytix.com

0 www.innocentive.com

0 tunedit.org

0 Tips for winning

0 [http://www.allanalytics.com/author.asp?doc\\_id=268513](http://www.allanalytics.com/author.asp?doc_id=268513)

# Mentors & Community

0 Google+

0 Data Science

0 Statistics and R

0 Artificial Intelligence

0 Machine Learning

0 Python

0 MongoDB

0 LinkedIN

0 Twitter

0 IRC

0 People within your  
company

# Conferences

0 Strata

0 strataconf.com

0 PyData

0 pydata.org

0 PyCon

0 us.pycon.org

0 Big Data Summit KC

0 BigDataSummitKC.org

0 Investigative Reporters & Editors Conference

0 www.ire.org/conferences/

# Introduction to Data Science Questions?

**Bryan Nehl**

@k0emt

dbBear.com

*(Links to Twitter, blog, GitHub, gists)*