# A Contextual introduction to Data Science

Bryan Nehl

@k0emt

Kick starting your own exploration!

# Skills

O Data Journal – Engineer's Notebook

  O http://soloso.blogspot.com/2014/05/engineers-notebook.html

O Regular Expressions, Markdown

O Analysis (math/stats is part of this!)

O LINUX/*NIX/ OS X/macOS

  O putty

O Version Control

  O http://git-scm.com/

O Messaging

  O http://rabbitmq.com

  O http://zeromq.org

# Virtualization/Containers

- Why?
  - Time
    - Pre-built images
  - Cost
  - On Demand

- How/Where?
  - Microsoft Azure
    - Data Science VM
  - Amazon Elastic Cloud
  - Google Compute
  - Your Own Machine
    - Oracle VirtualBox -
    www.virtualbox.org
  - Docker
    - hub.docker.com

# Common Virtual Machines (VMs)

**LAMP/WAMP**
- Linux/Windows
- Apache
- MySQL
- PHP

**MEAN**
- MongoDB
- Express
- Angular
- Node.JS
- meanjs.org
- meteor.com

mongoDB

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 by Google

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

node JS

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

# Azure Data Science VMs

- *Windows* Based VM
- Microsoft R Server Developer Edition
- Anaconda Python
- Jupyter notebooks
  - Python & R
- Visual Studio CE
  - Python & R Tools
- Power BI desktop
- SQL Server Express
- Machine Learning Tools

- *Linux* Based VM
- Microsoft R Open
- Anaconda Python
- Jupyter notebooks
  - Python & R
- Postgres Database
- Azure Tools
- Machine Learning Tools

# Journal

O Linux, Linux VM, mac     O Workstation

# Sourcing the data

o Locate it
  o Provided
  o Search for it
    o Manually
    o Automated
  o Networking

o Get it
  o ftp
  o Scraping
  o Database
  o Web services
o Work with it

# Some data

- data.gov
- data.mo.gov
- data.kcmo.gov
- data.gov.uk
- data.worldbank.org
- aws.amazon.com/datasets
- gutenberg.org
- https://gist.github.com/k0emt/63f19f828561c074f119
- soloso.blogspot.com/2011/07/getting-enron-mail-database-into.html

# Journal

- Sourcing Data
  - XML file layout example

# NEVER

*trust*

the data!

# Data Analysis ???

- What do I expect to see?
- What are the field types?
- Does the field type change?
- What are the range of values?
- How frequently do those values occur?
  - *Can I get a graph please?*
- Are there nulls?

- How big is my sample set?
  - *Is it significant?*
- How big do I expect the real data to be?
- Are there holes in the data?
- What constitutes a *good* record?
- Where are the trends/clusters in the data?

# Journal

O Upfront Analysis
   O File layout
   O File description

# Extract, Transform & Load
## Data Formats

- XLS
- CSV
- Text
  - Delimited
  - Fixed format
- JSON – json.org

- XML & HTML
- Mail files
- SQL scripted INSERTs
- PDF

Character sets
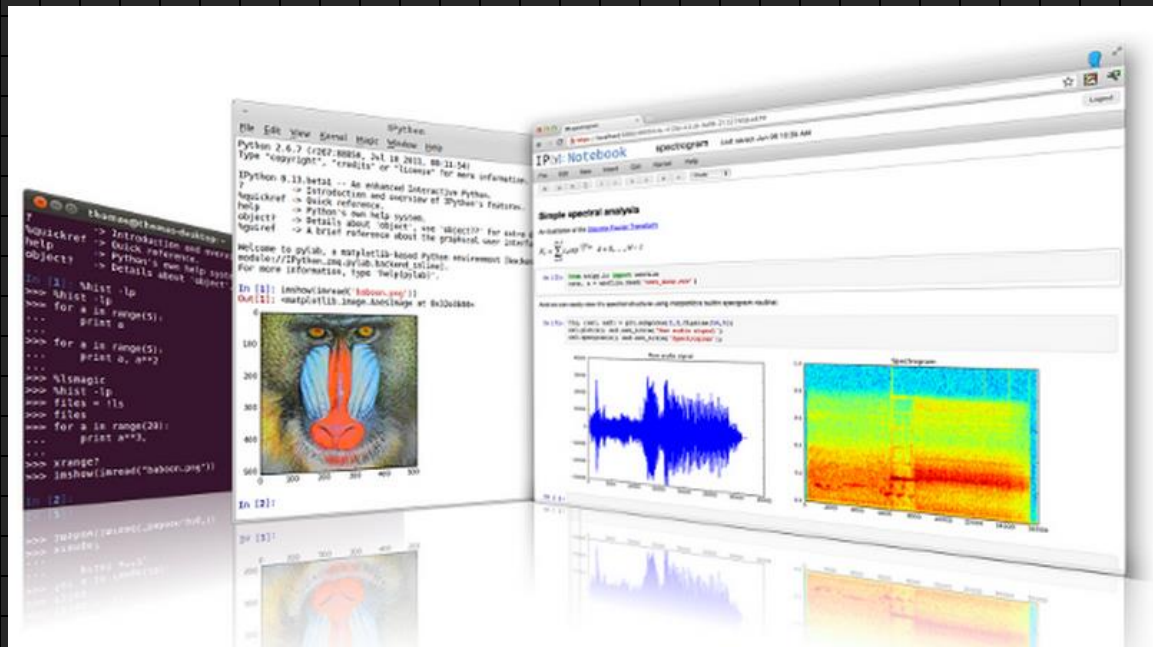https://docs.python.org/2/howto/unicode.html

# Languages and Libraries Extract *&Transform*

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

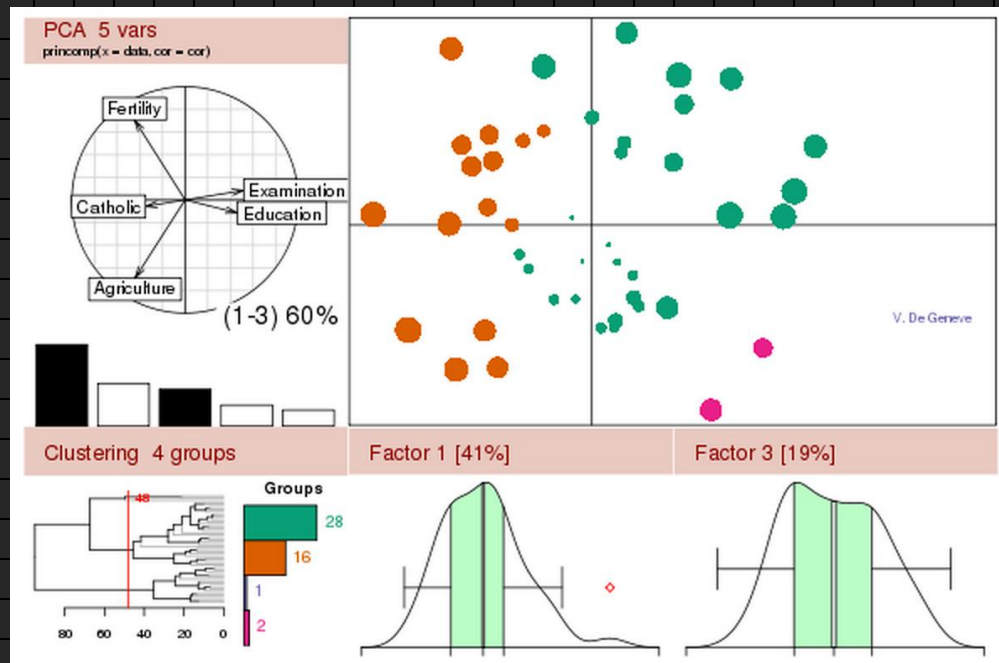# Languages

- Python – python.org
- IPython – ipython.org
- http://nbviewer.jupyter.org/

# Languages

O R - www.r-project.org

O R Studio - www.rstudio.com

# Libraries for Excel

O It is everywhere

O Python Libraries:

    O xlrd

    O XlsxWriter

O Apache Project – Office Open XML file formats

    O http://poi.apache.org/

      O Excel

      O Word

      O PowerPoint

# Libraries

- SciPy
  - scipy.org
- NumPy
  - numpy.org
- Pandas - Python data analysis library
  - pandas.pydata.org

# Libraries



O lxml
  O `lxml.de`
O pymongo
  O `pypi.python.org/pypi/pymongo/`
O pika – AMQP
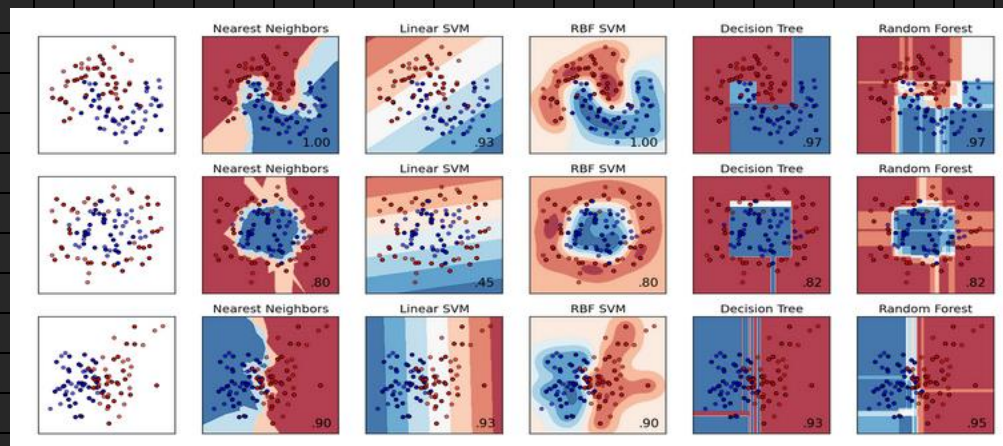  O `pypi.python.org/pypi/pika`
O nose – unit test framework extension
  O `nose.readthedocs.org`

# scikit-learn.org

O Machine Learning
  O Clustering
  O Classification
O 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://www.continuum.io/why-anaconda

0 wakari.io web based Python data analysis

# Databases

Choose the right one(s) for the job!

Polyglot Persistence
http://martinfowler.com/bliki/PolyglotPersistence.html

# Relational - SQL

O MySQL
  O open source
O Oracle
O Microsoft SQL Server
O Express Editions
O Microsoft Access


O ODBC / JDBC

# NOSQL

O Definition

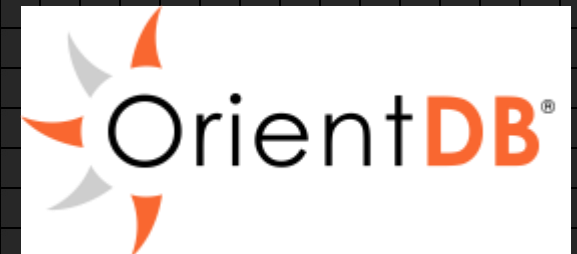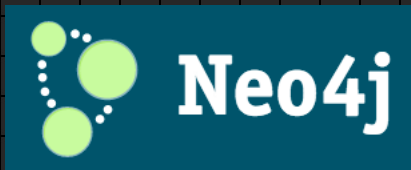O MongoDB – JSON/BSON documents
  O mongodb.org
O neo4j – graph
  O neo4j.org
O OrientDB – document & graph
  O orientdb.com
O PostgreSQL – object-relational
  O postgresql.org

O Distributed framework for processing large datasets

O MongoDB and other databases can be used to feed it

O MapReduce

O hadoop.apache.org



O In Memory MapReduce

O spark.apache.org

drill.apache.org

# Journal

- DOC ETL
  - sample data file
  - Counties
  - Python Program for translating into JSON
  - error file
  - import
- Analysis of available data
  - the top 5

# Business Context

I have data.  Now what?

# Numbers need context

| | | |
|---|---|---|
| Visitors | 1M | Last Year 2M |
| Page Views | 5.2M | Last Year 7.2M |
| 72% | Conversion Rate | |
| 42 | Customer average age | |
| 1 | Top Referrer.com | |

# Analysis

**Techniques**

o Adjacency Matrix

o pivot and fold operations on tables

o hexagonal binning

o confusion matrix

o predictive modeling fundamentals

o machine learning

o The work of John Tukey (Statistics)

  o wikipedia.org/wiki/John_Tukey

# Looking at numbers

| | | | | | |
|---|---|---|---|---|---|
| 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 |

Conditional Formatting – Color Scales

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

O Is there some insight I can bring?

O Can I generalize from this data? (global)

O Can I ascertain local area insights?

O Are there natural partitions in the data?

    O Gender, race, age, location?

O Is there some business pain I can relieve?

O Can I enhance an existing data set?

O Can I bring in the data product with a shorter cycle time?

# The Science Part

o Ask a question

o Form a hypothesis

o Do the research

# Journal

O Questions

O Data Work to Answer the questions

   O population buckets

# Visualization

Use your pixels!

# HTML Tools & Libraries

O HTML5 / CSS3

O Javascript

O **D3 – d3js.org**

O HTML 5 canvas charts

O chartjs.org

O canvasjs.com

# Google Tools & Libraries
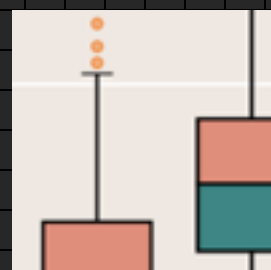
o Google Charts

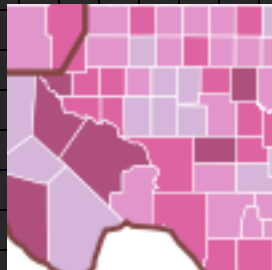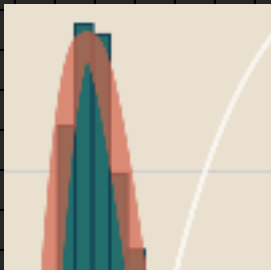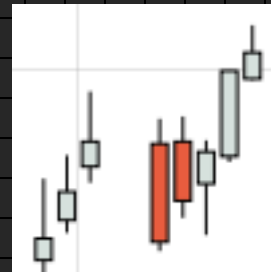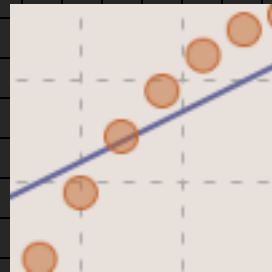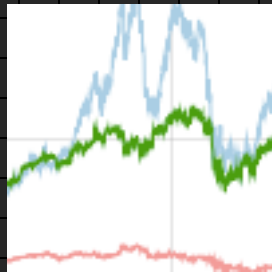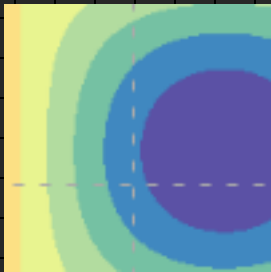   o developers.google.com/chart/



o Google Fusion Tables

   o Now integrated with Google Drive
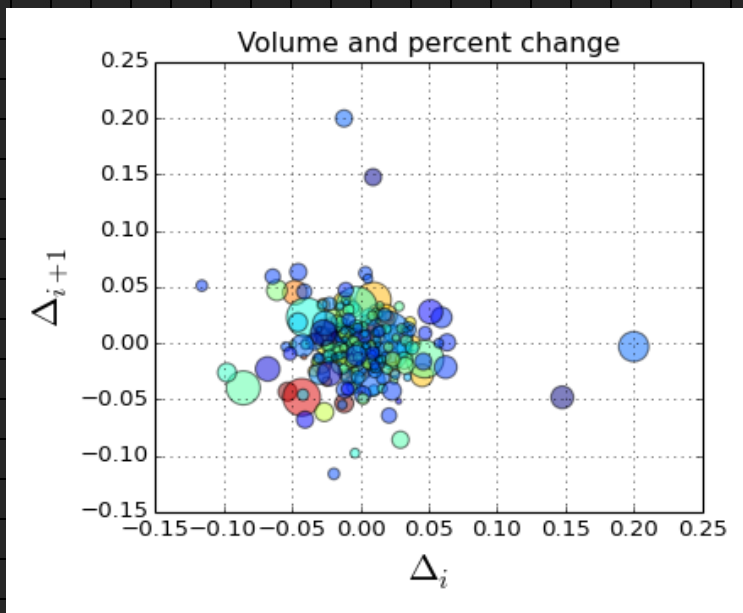
# Python Tools & Libraries

*o* bokeh - bokeh.pydata.org

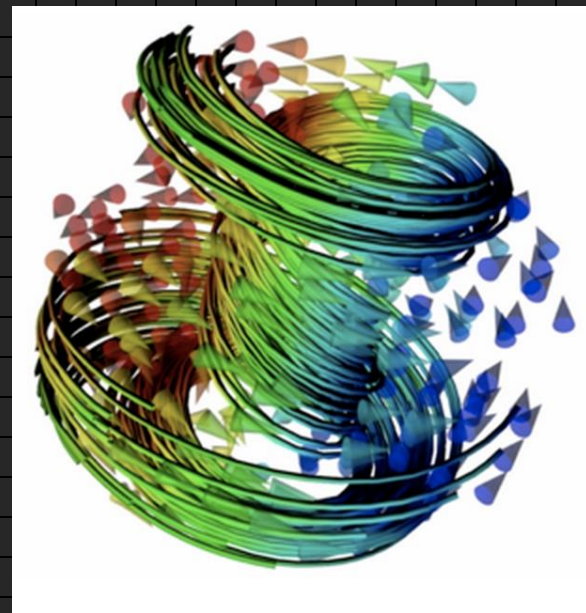# Python Tools & Libraries

○ Matplotlib
  ○ matplotlib.org



○ Mayavi 2
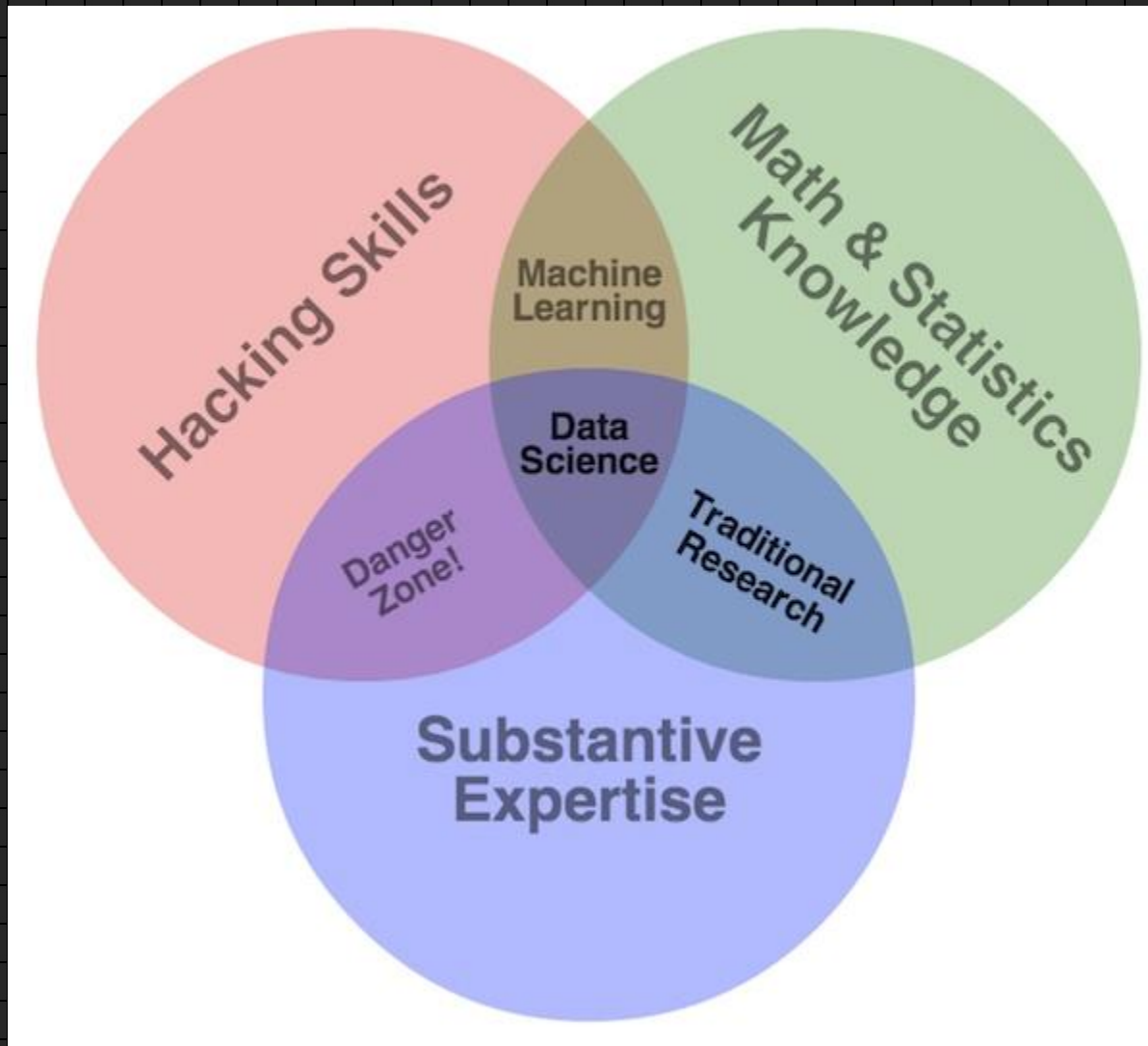  ○ code.enthought.com/projects/mayavi/

# Journal

O Visualization

   O Excel

   O R

     O Box Plots - base

     O Violin Plots – ggplot2

O Anaconda & bokeh

## Don't jump to conclusions!

*What now?*

# Drew Conway's Diagram

# Data Science Venn Diagram v2.0

Data Science

Computer Science

Machine Learning

Math and Statistics

Unicorn

Traditional Software

Traditional Research

Subject Matter Expertise

# Teams

# Decide your direction

○Personal Tech Radar

http://nealford.com/memeagora/2013/05/28/build_your_own_technology_radar.html

# Conferences

o Strata – `strataconf.com`

o Open Data Science Conference – `odsc.com`

o PyData – `pydata.org`

o PyCon – `us.pycon.org`

o Big Data Summit KC – `BigDataSummitKC.org`

o Investigative Reporters & Editors Conference – `www.ire.org/conferences/`

# Training

O Tutorials and sample files that come with software.

O Local courses

O Online Education from vendors

   O MongoDB University

      O university.mongodb.com

O Other online education

   O Coursera – coursera.org

   O iTunes University (iTunes U)

   O O'Reilly Safari, books, videos and free publications

   O oreilly.com/data/free

   O YouTube

   O Open Source Data Science Masters

      O datasciencemasters.org

# Mentors & Community

O Google+
  O Data Science
  O Statistics and R
  O Artificial Intelligence
  O Machine Learning
  O Python
  O MongoDB

O LinkedIN
O Twitter
O IRC
O People within your company
O BecomingADataScientist.com
O Reddit /r/datascience
O Datascience.stackexchange.com

# Contests

O kaggle.com

O www.kdnuggets.com/competitions/

O www.crowdanalytix.com

O www.innocentive.com

O tunedit.org

O drivendata.org/competitions/


O Tips for winning
   O http://www.allanalytics.com/author.asp?doc_id=268513

# Experiment

O Set up a development environment

O Create a Virtual Machine

O Spin up containers

O Try out stuff

   O Work related

   O Something you are passionate about

O Share your experiences

   O blog, tweet, present

   O GitHub and Gists

*Enjoy your journey!*

Bryan Nehl – **@k0emt** – dbBear.com

https://github.com/k0emt/Presentations

https://github.com/k0emt/**corrections**