

Data Science with Python

Phil Mui
@philmui

Today

- Hello
- Data Science: what is it?
- Course Overview
- Python: codeacademy
- Jupyter / ipython
- Git / VirtualEnv
- Check-in

Today

- Hello
- **Data Science: what is it?**
- Course Overview
- Python: codeacademy
- Jupyter / ipython
- Git / VirtualEnv
- Check-in

What is Data Science?

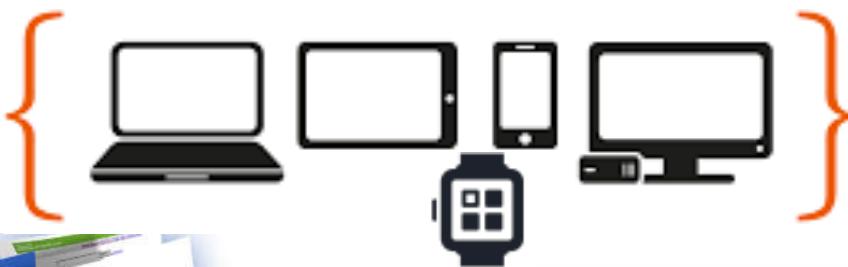
- What words come to mind when you think of Data Science?
- What experience do you have with Data Science?
- Why are you taking an Introduction to Data Science Class?

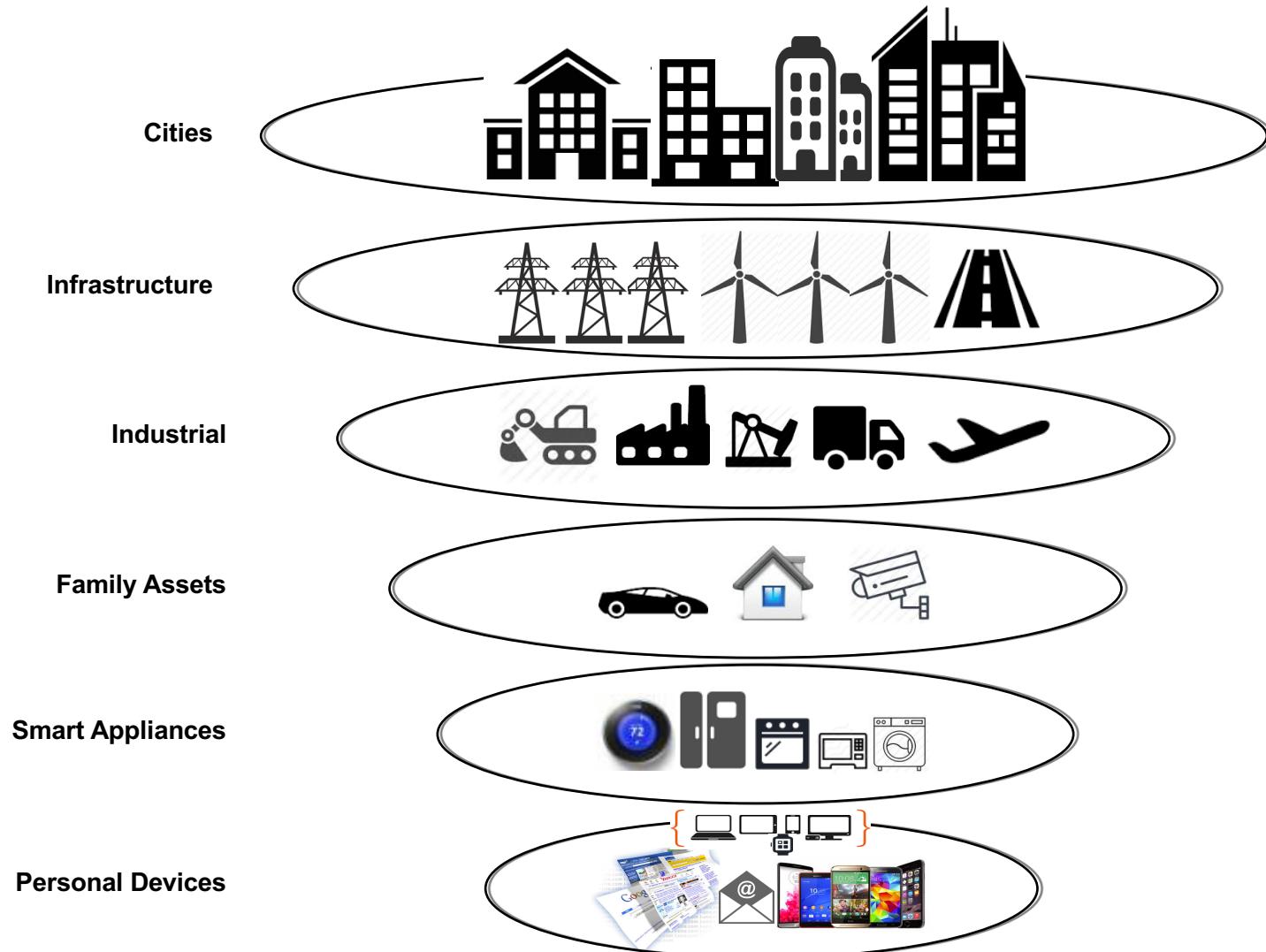
“Data”











HBR.ORG

Harvard Business Review

OCTOBER 2012
REPRINT R1310C

SPOTLIGHT ON BIG DATA

Big Data: The Management Revolution

Exploiting vast new flows of information
can radically improve your company's
performance. But first you'll have to
change your decision-making culture.
by Andrew McAfee and Erik Brynjolfsson

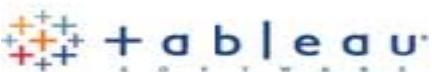
Data Scientist:

The Sexiest Job of the 21st Century



REVO^{LUTION}
ANALYTICS



ClearStory^{DATA}
GoodData
The Tableau logo consists of a grid of colored dots (red, green, blue) followed by the word 'tableau' in a lowercase sans-serif font.

Quantivo

alteryx

webtrends



Google Analytics



mixpanel

Localytics

Spark



ASTER



cloudera

Greenplum

NETEZZA

VERTICA

ORACLE

SAP

IBM

TERADATA

eBureau

AnalyticsIQ Inc.

POLK

IXI™ SERVICES

causa
nielsen

Semcasting

What is Data Science?

- “How Companies Learn Your Secrets” *NYT*, by Charles Duhigg, February 16, 2012
- <http://bit.ly/targetdata>



<http://bit.ly/targetdata>

What did Target Do?

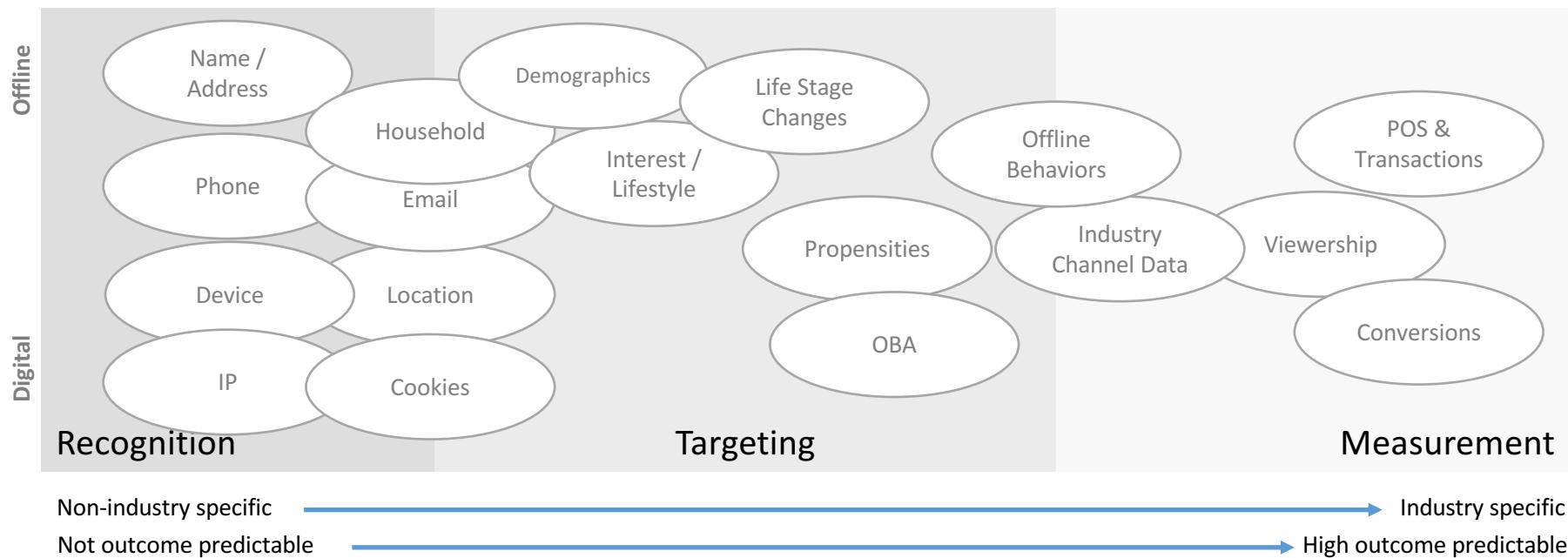
- Mining of data on shopping patterns
 - Specific products purchased
 - Combination of products purchased
 - Combined with demographic and other data
- Psychology and neuroscience
 - Habits:
 - Cue-routine-reward
 - When are habits open to change?

Lessons from Target

- Yes, Data Science is about mining data
- There are deeper theoretical issues involved in understanding what you find
- Left out of that long article are most of the critical steps that precede the analysis
- In short, Data Science > data mining

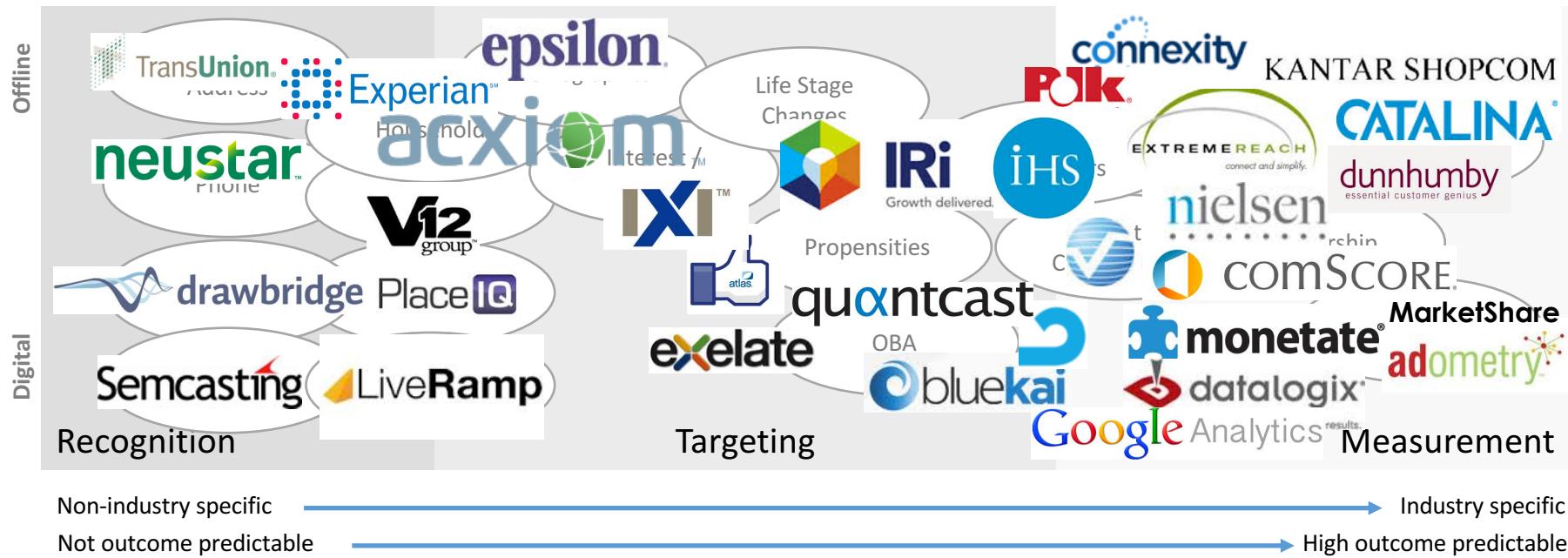
Consumer data landscape

A framework for understanding consumer data



Consumer data landscape

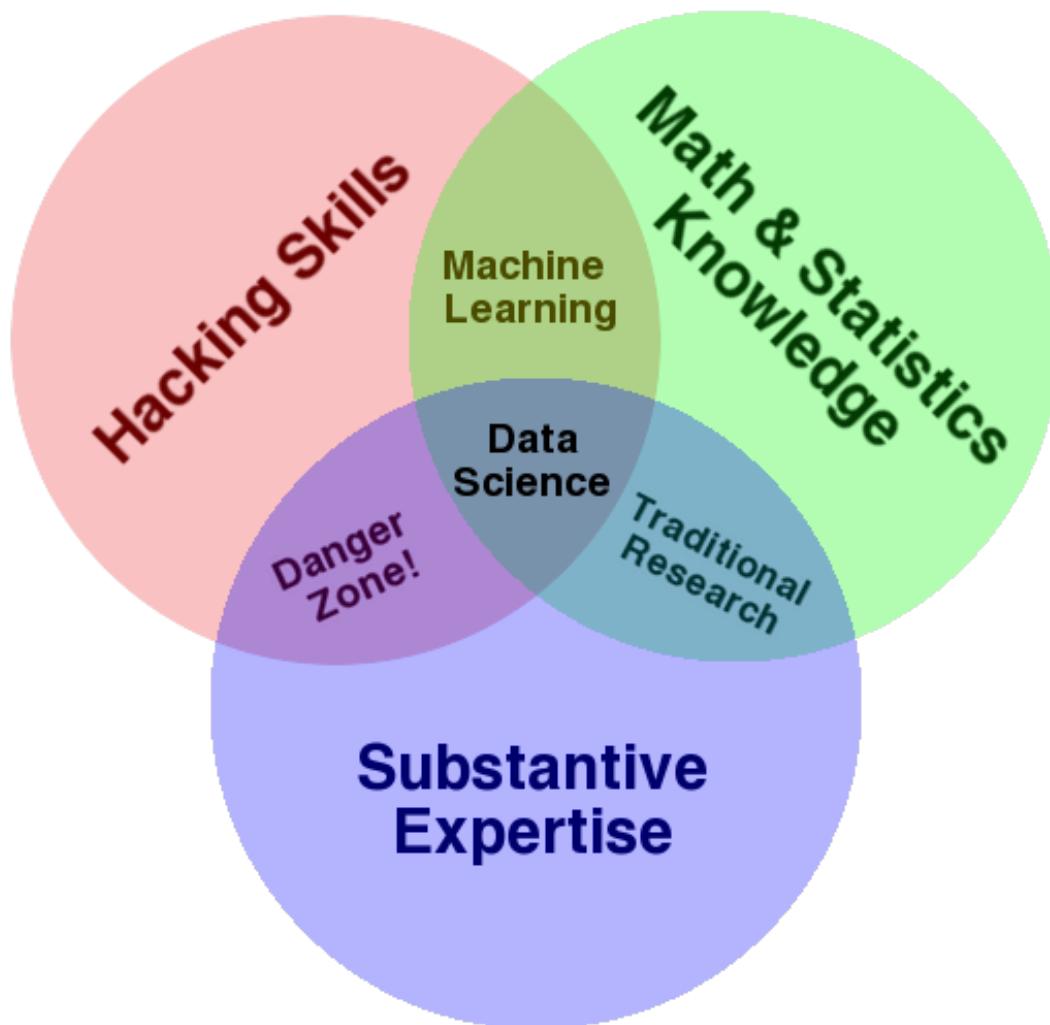
A framework for understanding consumer data



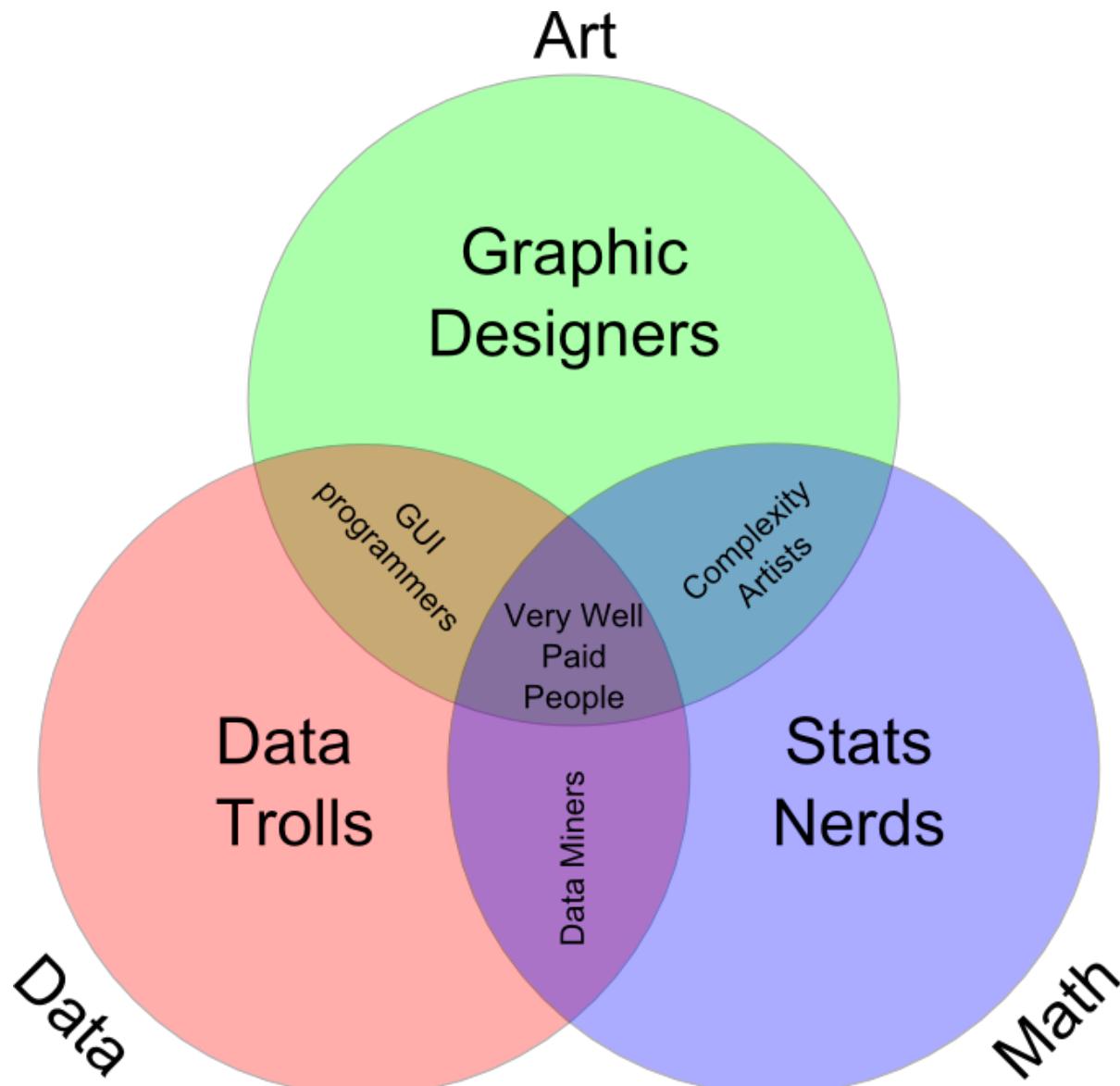
Definition of Data Science

- There are many, but most say data science is:
 - Broad – broader than any one existing discipline
 - Interdisciplinary: Computer Science, Statistics, Information Science, databases, mathematics
 - Applied focus on extracting knowledge from data to inform decision making.
 - Focuses on the skills needed to collect, manage, store, distribute, analyze, visualize, and reuse data.
- There are many visual representations of Data Science

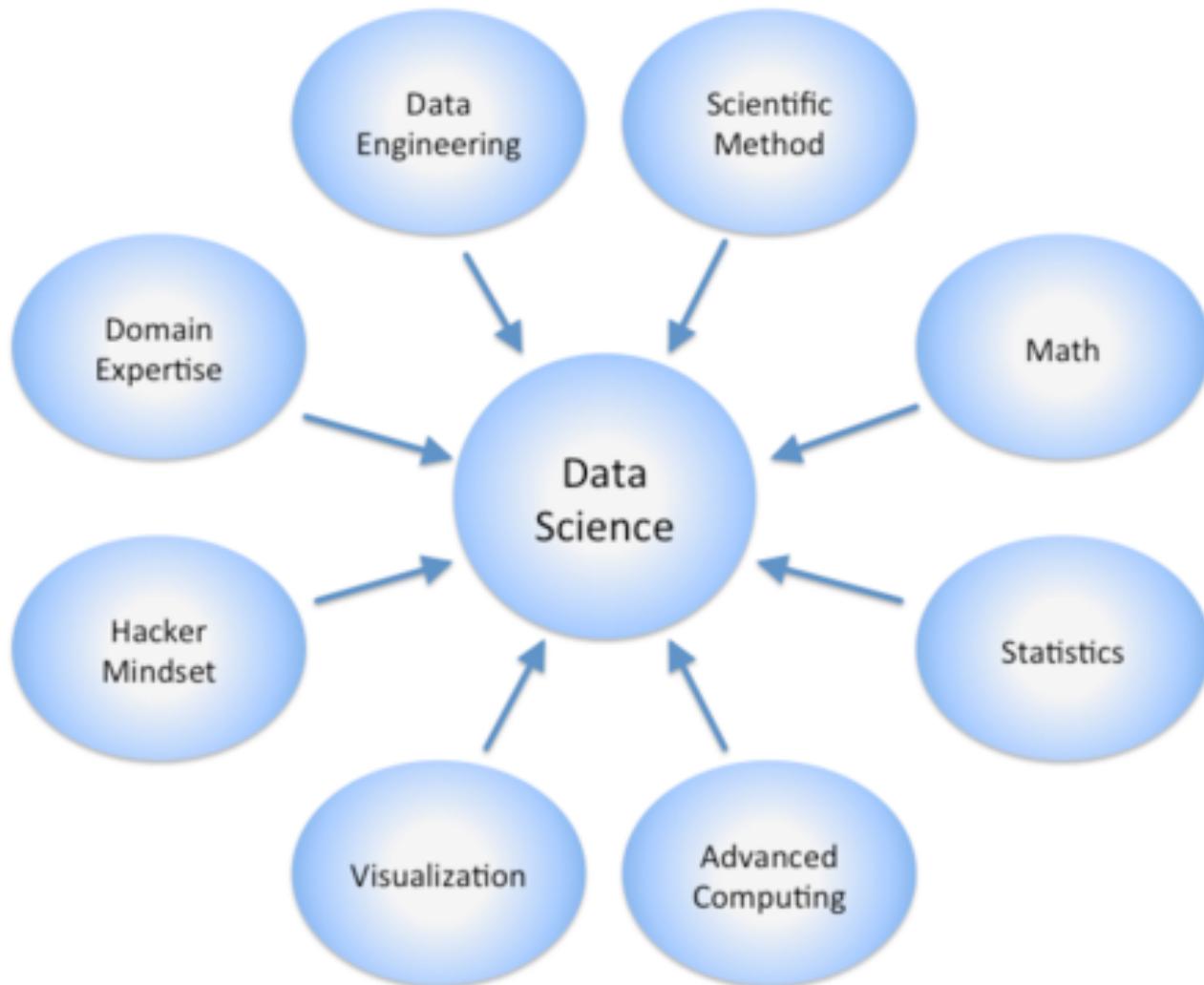
Some definitions link computational, statistical, and substantive expertise.



Other definitions focus more on technical skills alone.



Still other definitions are so broad as to include nearly everything.

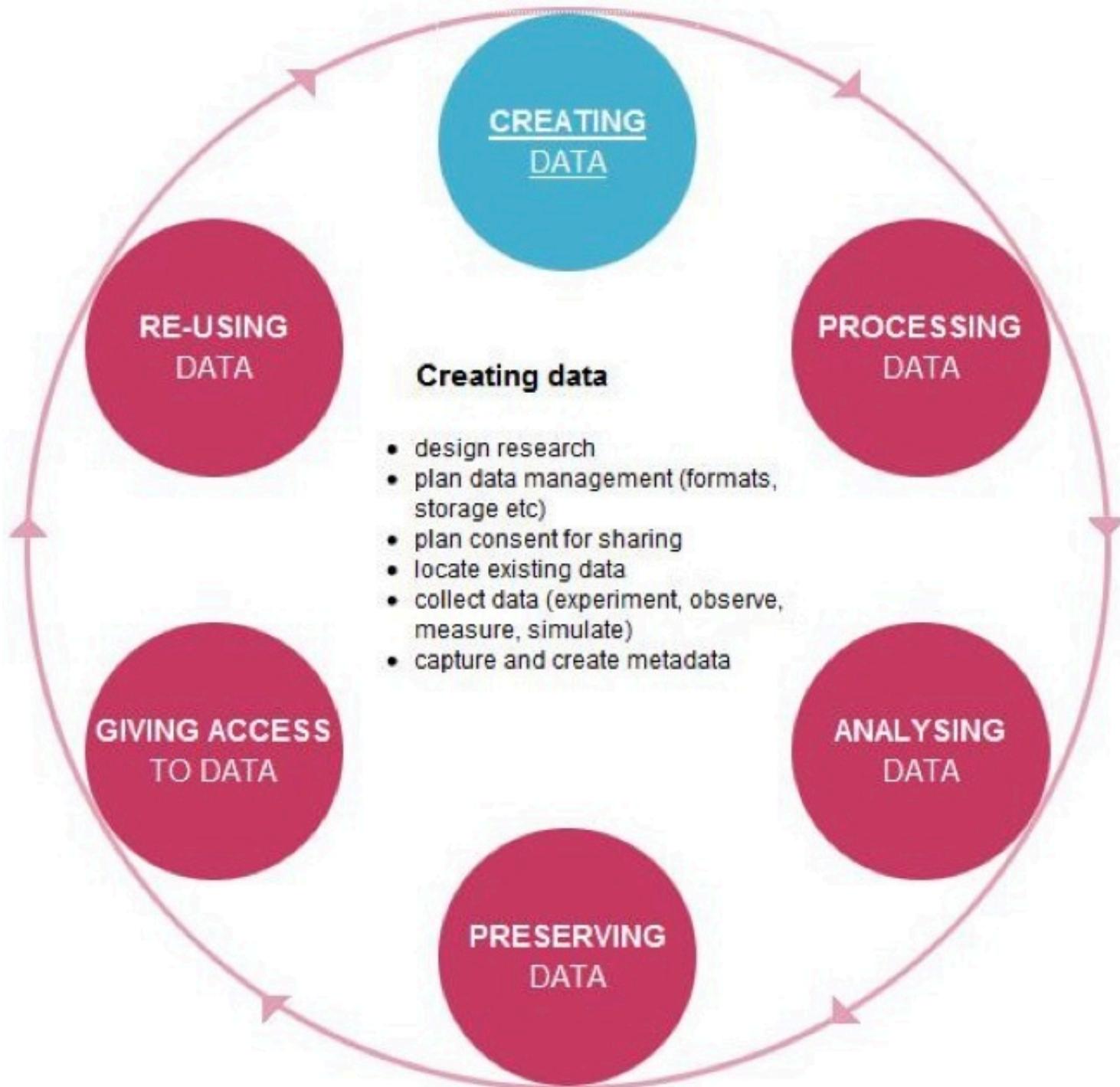


There are many “Word Cloud” representations of Data Science as well.



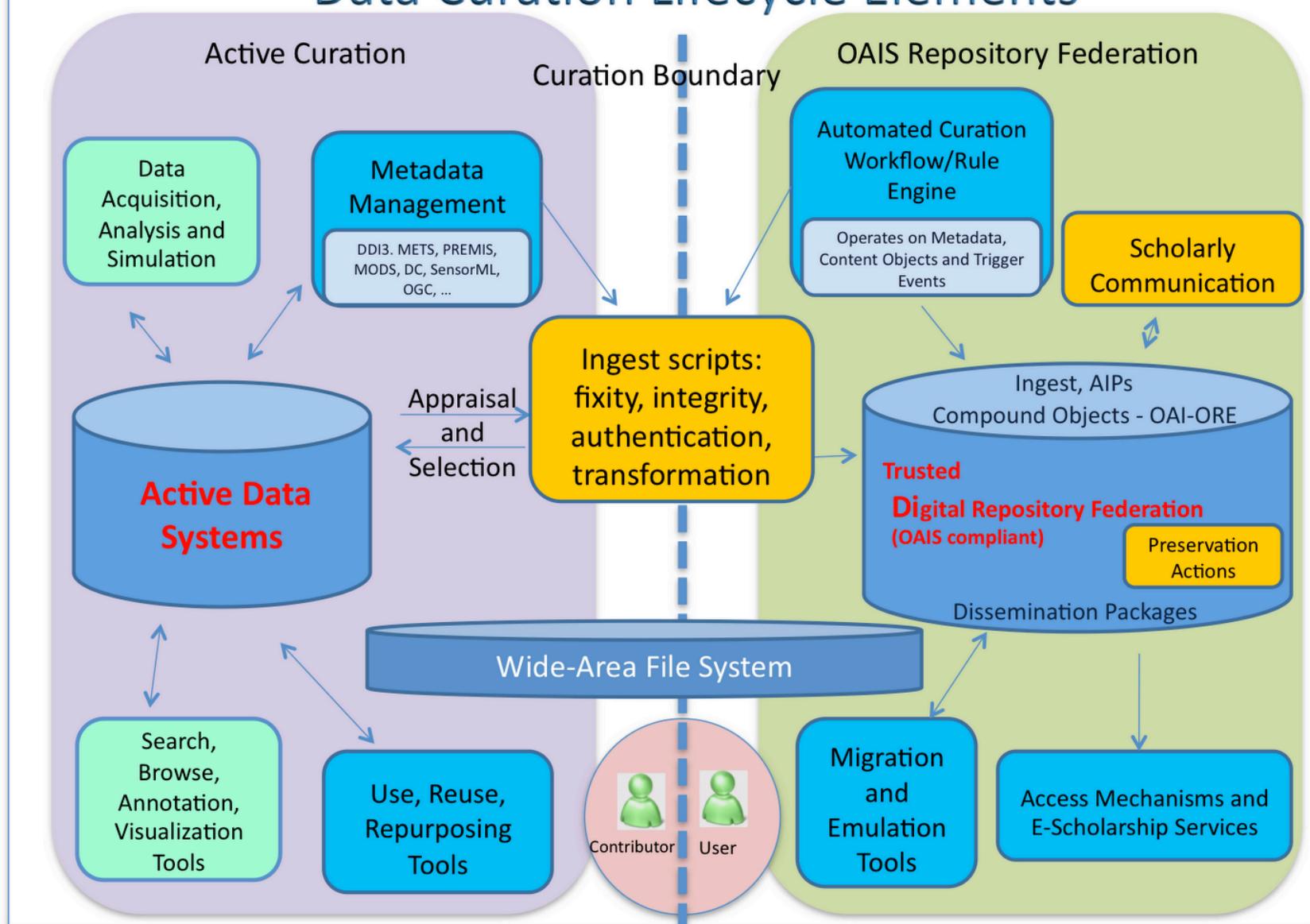








Data Curation Lifecycle Elements



What is Missing?

- Most definitions of data science underplay or leave out discussions of:
 - Substantive theory
 - Metadata
 - Privacy and Ethics

Today

- Hello
- Data Science: what is it?
- **Git / VirtualEnv**
- Course Overview
- Python: codeacademy
- Jupyter / ipython
- Check-in

Course github site:

<https://github.com/philmui/datascience2016fall>

```
git clone https://github.com/philmui/datascience2016fall
```

Getting git

Mac:

- brew install git

PC:

- <https://git-scm.com/downloads>

Today

- Hello
- Data Science: what is it?
- Git / VirtualEnv
- **Course Overview**
- Python: codeacademy
- Jupyter / ipython
- Check-in

docs/Data.Science.Python.Syllabus.ipynb

Today

- Hello
- Data Science: what is it?
- Git / VirtualEnv
- Course Overview
- **Python: codeacademy**
- Jupyter / ipython
- Check-in

Getting on the same page on Python



<https://www.codecademy.com/learn/python>

Goal: finish by October 15th

- We have no in-person class on October 15th
- You should use that day to complete the ENTIRE CodeAcademy Python course
- Submission to Camino a screenshot of your completed Python course (no need to pay)

Today

- Hello
- Data Science: what is it?
- Git / VirtualEnv
- Course Overview
- Python: codeacademy
- **Jupyter / ipython**
- Check-in

lecture01/lecture01.intro.ipynb

Jupyter / IPython

The screenshot shows a Jupyter Notebook interface running in a web browser window titled "lecture01.jupyter". The browser's address bar displays the URL "localhost:8888/notebooks/datascience2016fall/lecture01.intro/lecture01.jupyter.ip...". The notebook title is "lecture01.jupyter (autosaved)". The toolbar includes File, Edit, View, Insert, Cell, Kernel, Help, and a Python 2 kernel selector. Below the toolbar is a toolbar with various icons for file operations like new, open, save, and cell execution.

A quick tour of Jupyter (aka IPython) Notebook

You can start this Jupyter (aka IPython) notebook from a terminal by running

```
jupyter notebook --pylab inline
```

The --pylab inline is for plotting graphs

First, we need to explain how to run cells. Try to run the cell below!

```
In [ ]: import pandas as pd  
print "Hi! This is a cell. Press the ▶ button above to run it"
```

You can also run a cell with Ctrl+Enter or Shift+Enter. Experiment a bit with that.

Today

- Hello
- Data Science: what is it?
- Git / VirtualEnv
- Course Overview
- Python: codeacademy
- Jupyter / ipython
- **Check-in**

Assignment 1

- Open Jupyter notebook: “lecture01.intro.ipynb”
- Examine the EU Revolving Loan data set:
“eu_revolving_loans.csv”
- Answer the questions:
 - Which EU Country has the largest revolving loan % increase from 2000 – 2016?
 - Which EU Country has the largest revolving loan % decrease from 2000 – 2016?
 - Qualitatively describe in 5 sentences or less why these 2 countries may be seen these big changes during this period.
- Submission: to “Assignment” section of course Camino
- Due date: 10/08/2016 (Sat) 11:59pm PT