

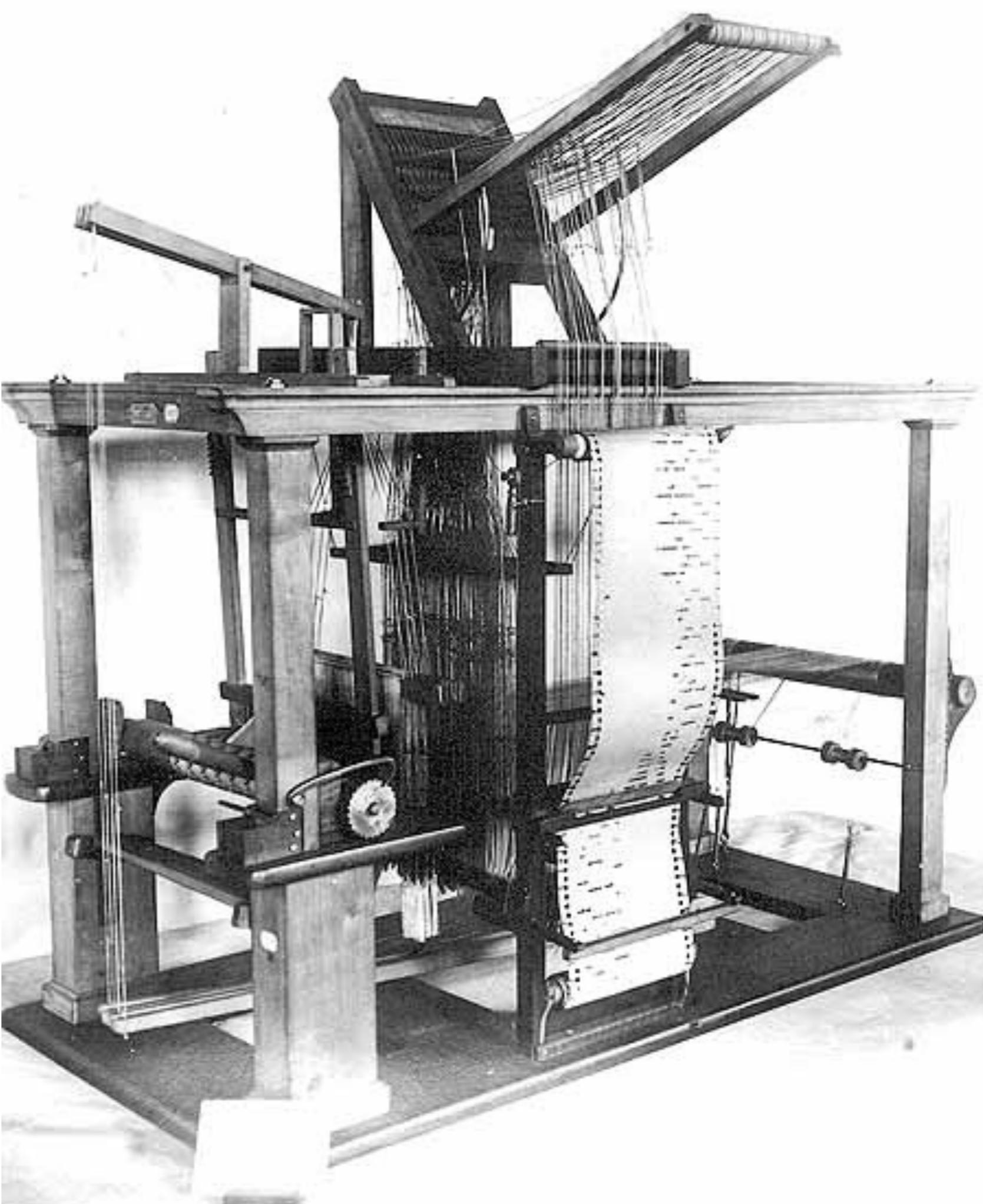
EDCT GE2550.

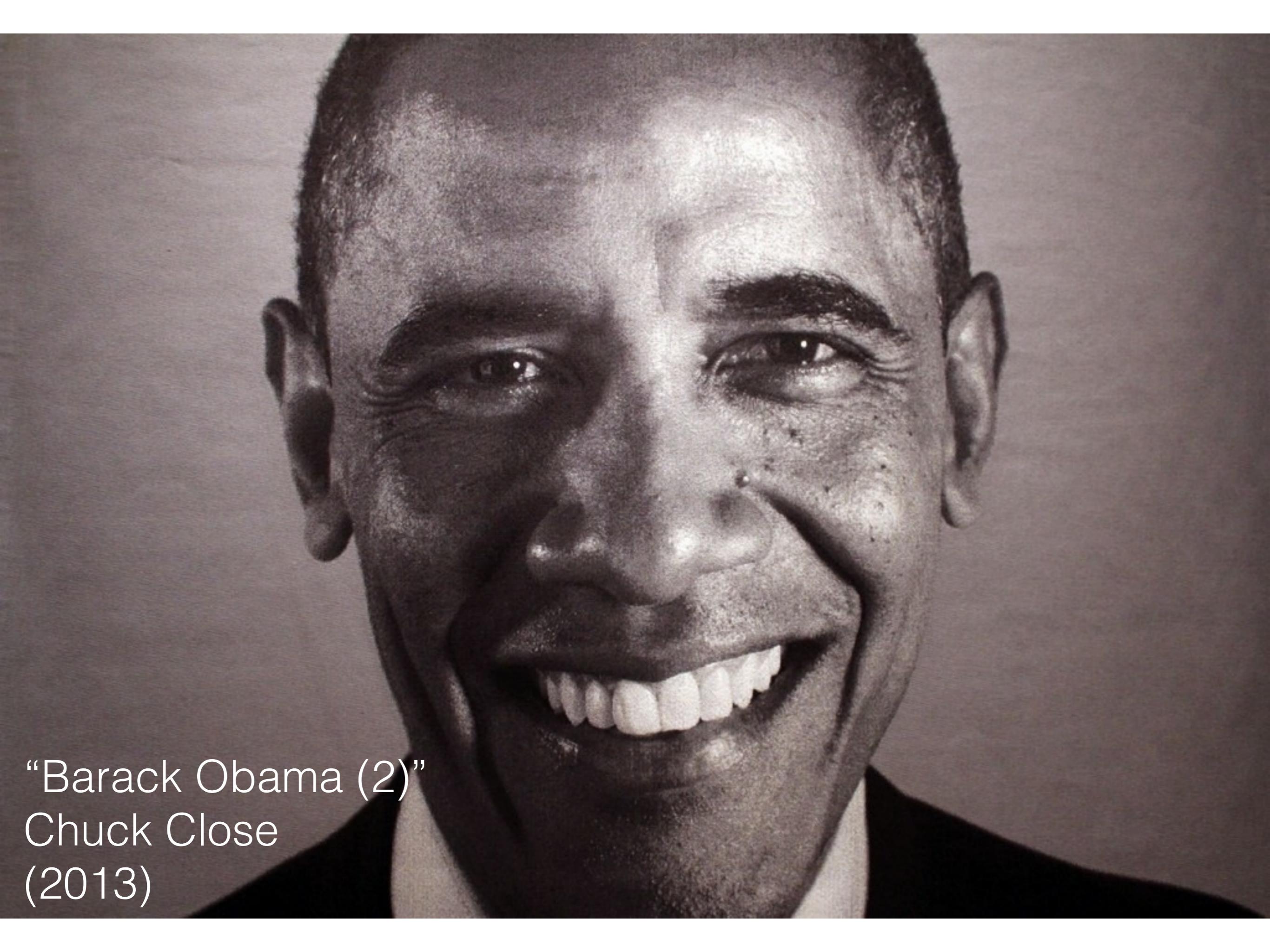
DATA SCIENCE

IN EDUCATION

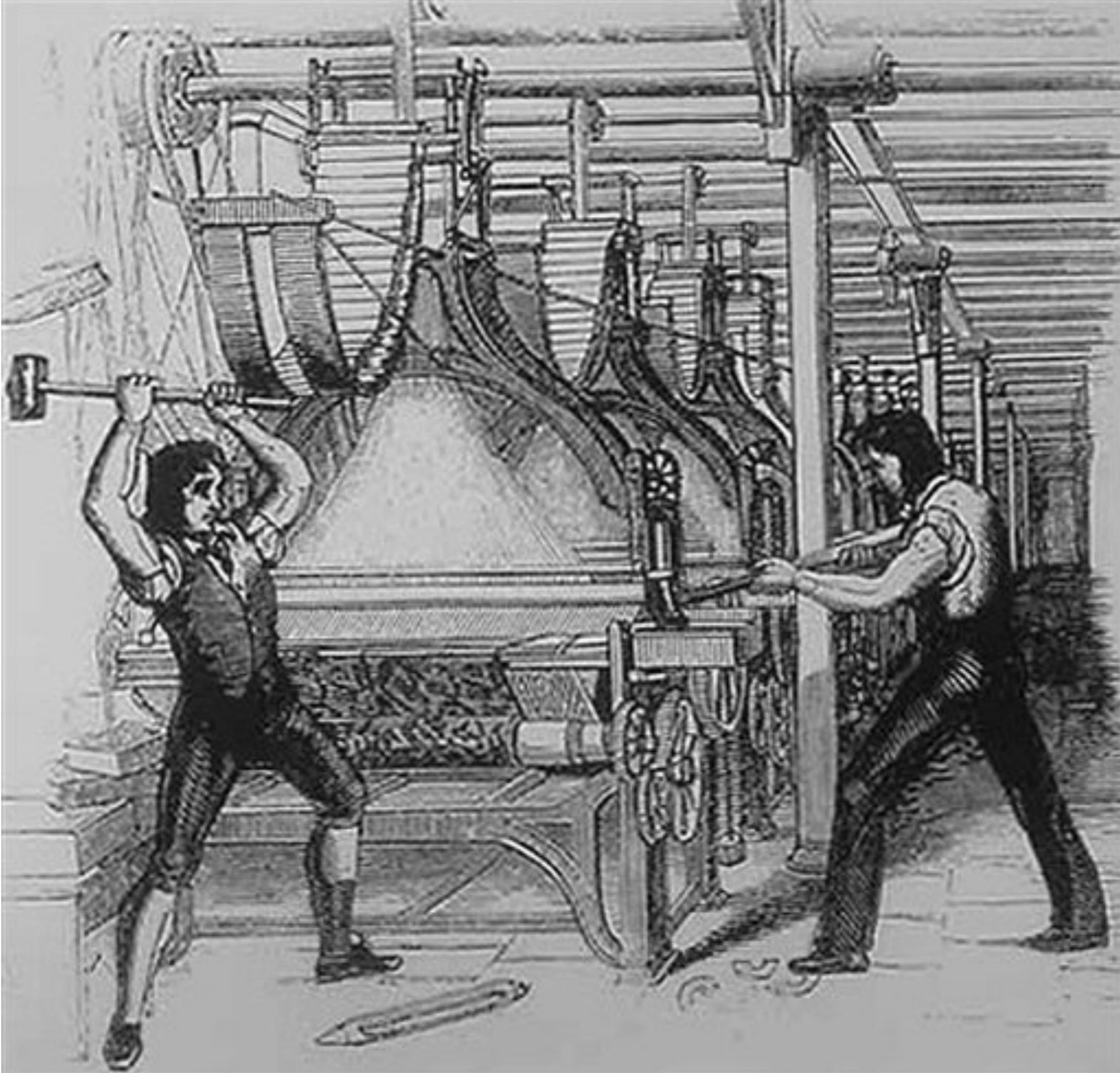
Big Data, Learning Analytics & The Information Age

Jacquard Loom  
(1801)





"Barack Obama (2)"  
Chuck Close  
(2013)



Luddites  
(1811-17)

Tweets

Media

Likes



leonie haimson @leoniehaimson

2d

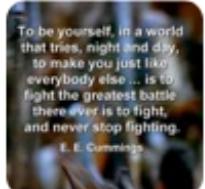
Treating kids like cattle  
@Parents4Privacy

Lori Lalama @TechEducator1

#DataMining #PrivacyAware  
@leoniehaimson @DianeRavitch  
[wired.com/2012/09/rfid-c...](http://wired.com/2012/09/rfid-c...)



Retweeted



Sheila Resseger @sheilaresseger

2d

@leoniehaimson I had no illusions 2  
the contrary! We r up against a  
relentless & callous assault on  
human dignity, preventing self-  
empowrment



Mark Zuckerberg

December 3, 2015 at 5:24pm · Palo Alto,  
CA ·

I want to thank you all for your heartwarming congratulations on Max's birth and on starting the [Chan Zuckerberg Initiative](#). This whole community has been so loving and supportive. If you're interested in following the philanthropy work we're doing with the Chan Zuckerberg Initiative, I encourage you to like the page here:

<https://www.facebook.com/chanzuckerberginitiative>

Since we announced this a couple days ago, many people have asked about what we're planning to focus on and how we're structuring our work.

Our initial focus areas are [personalized learning](#), curing disease, connecting people and building strong communities. We've already made many investments over the past five

## 1

## THE RAPID GROWTH OF GLOBAL DATA

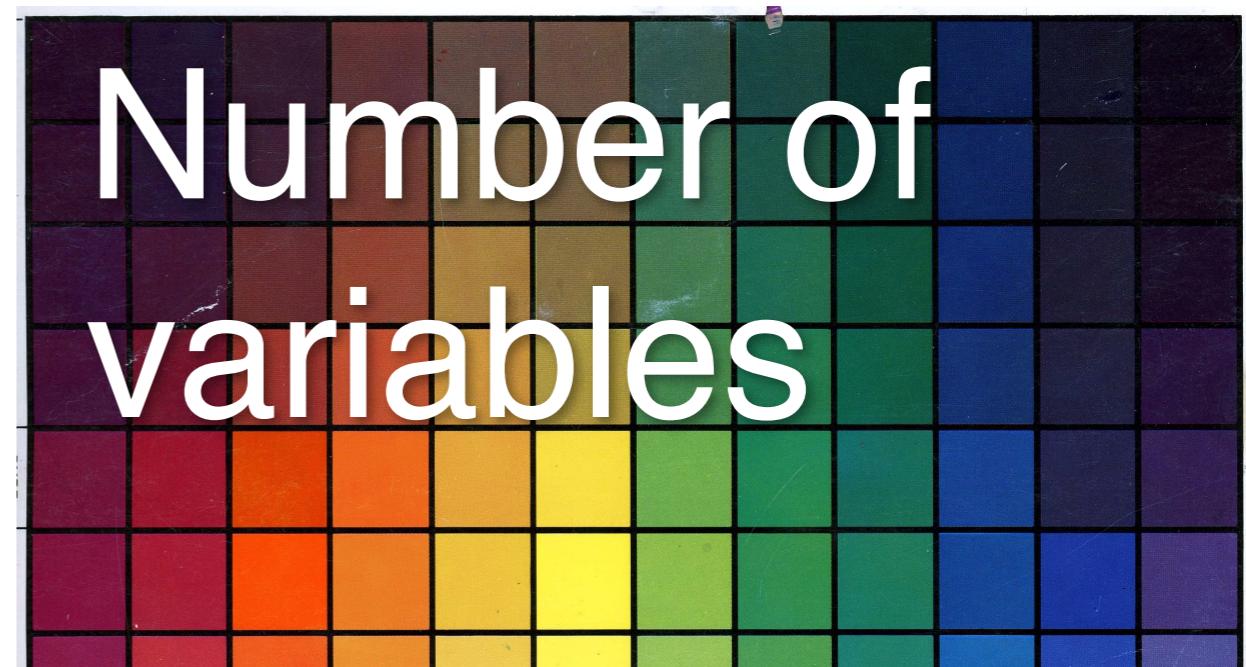
CSC

The production of data is expanding at an astonishing pace. Experts now point to a 4300% increase in annual data generation by 2020. Drivers include the switch from analog to digital technologies and the rapid increase in data generation by individuals and corporations alike.

**2020: MORE THAN 1/3  
OF THE DATA PRODUCED  
WILL LIVE IN OR PASS  
THROUGH THE CLOUD.**



# What's new?



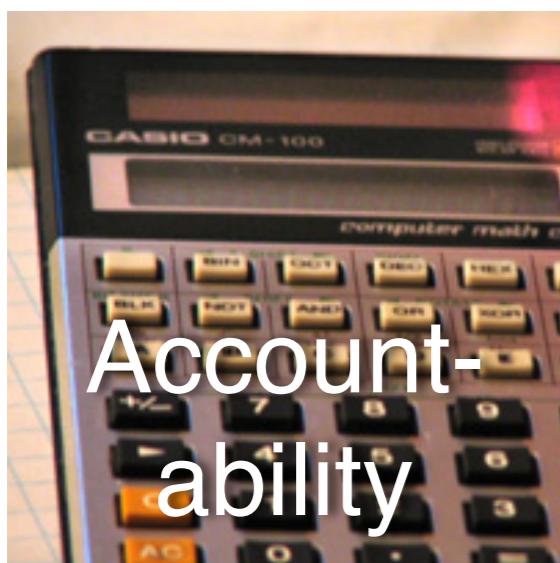
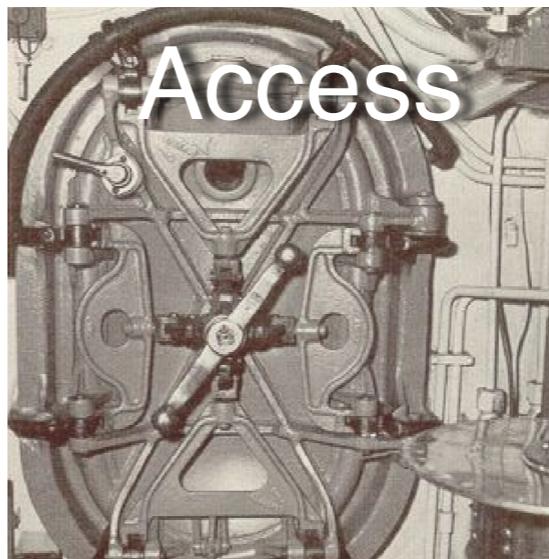
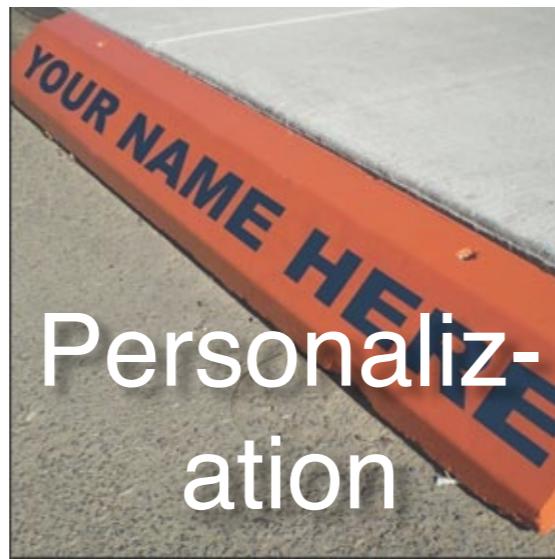
# MR. MESSY

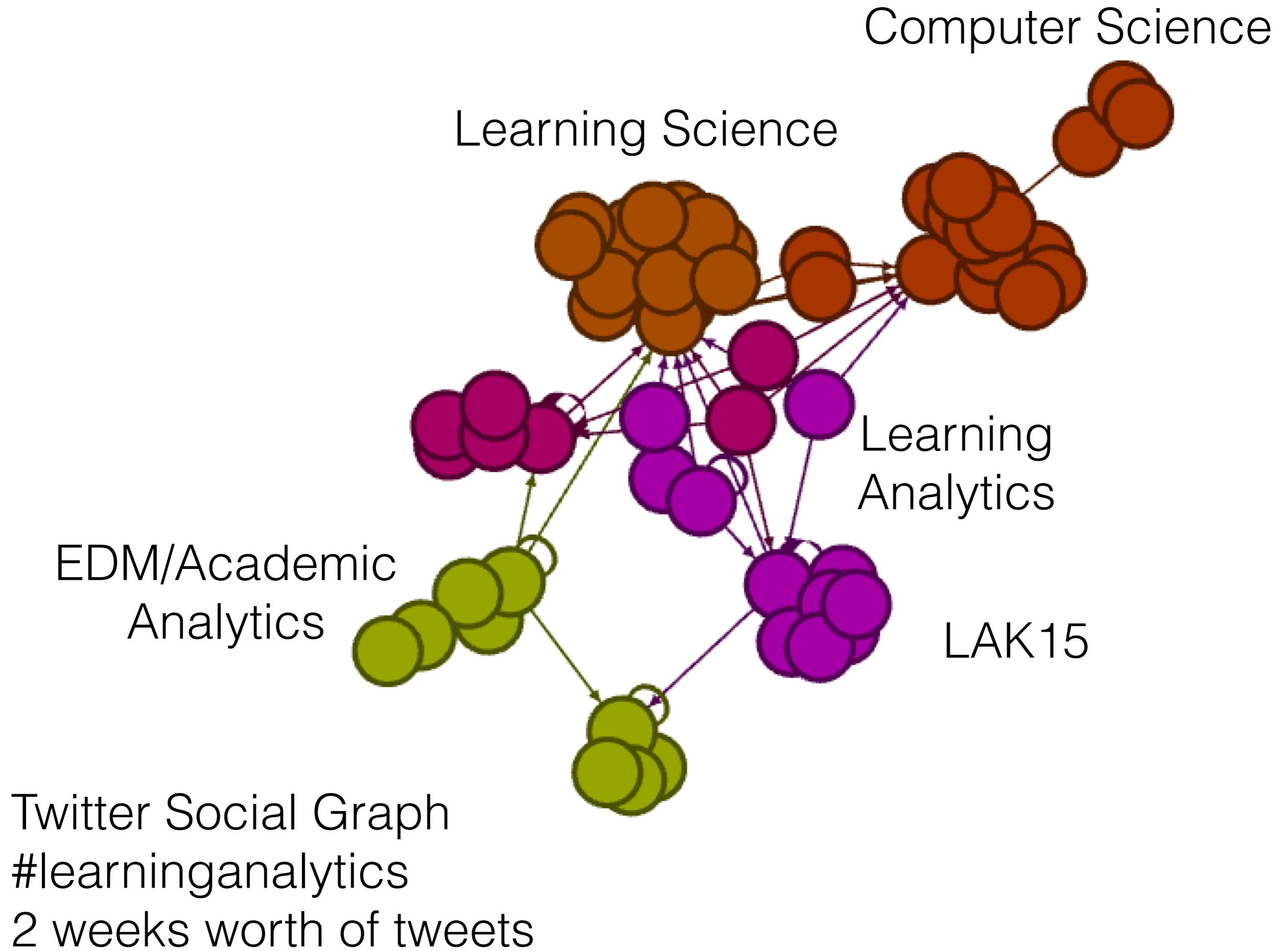
by Roger Hargreaves



- Volume
- Velocity
- Variety
- Veracity

# Possibilities





## Learning Analytics

- Grew from BI
- Early-mid 2000s in Canada
- Large domain umbrella - EDM, MOOCs, learning sciences, game design
- Systems oriented
- Methodologically broad (qual, quant, CS, stat)
- Journal: JAL
- Society: SOLAR

## Ed. Data Mining

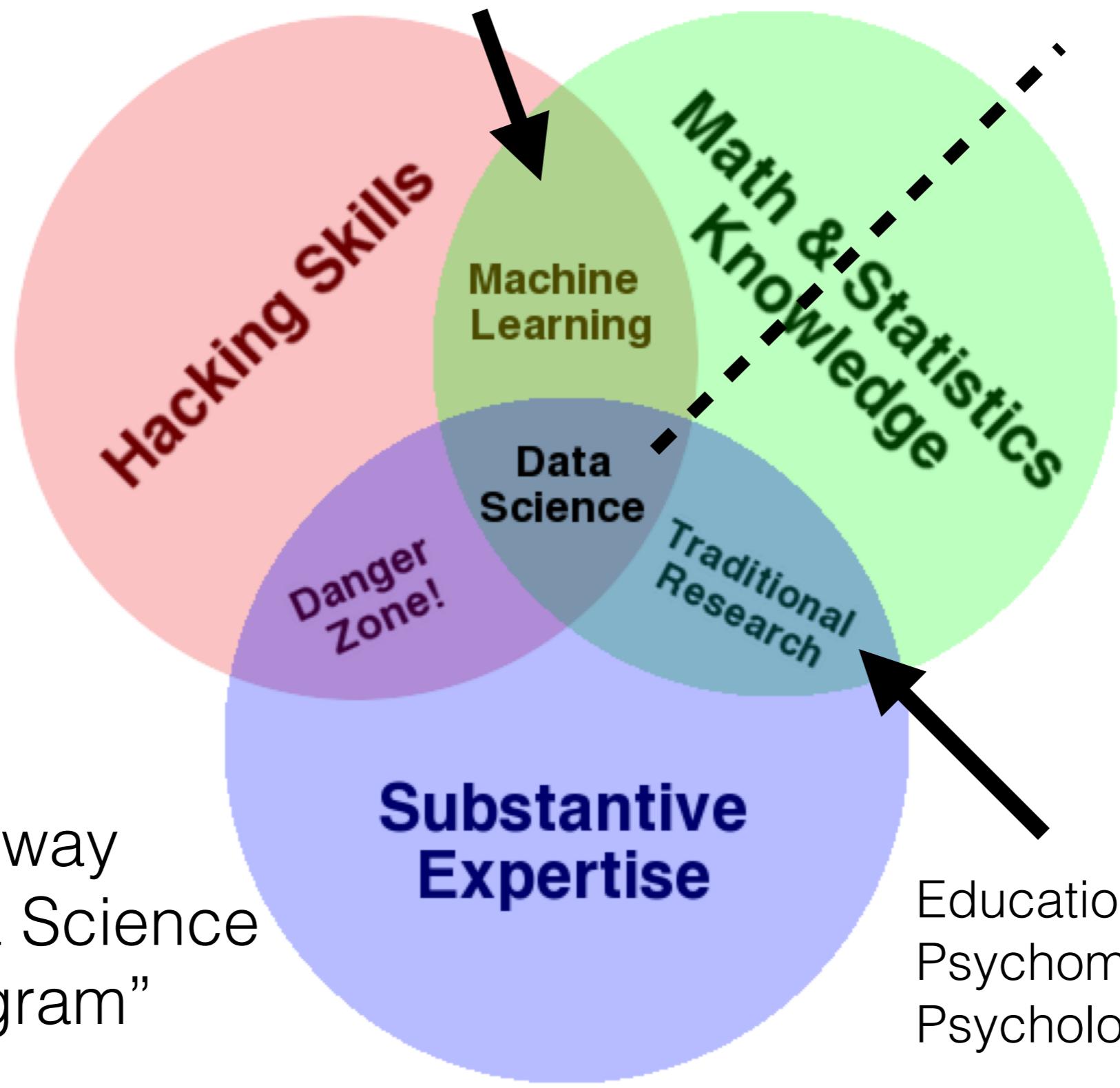
- Grew from Knowledge Mining or KDD
- Late 90s at CMU
- Intelligent tutoring, log-file
- Software oriented
- CS methods
- Journal: JEDM
- Society: IEDMS

“A Data Scientist is a statistician who lives in San Francisco”

–Sean Owen (Cloudera), 2014



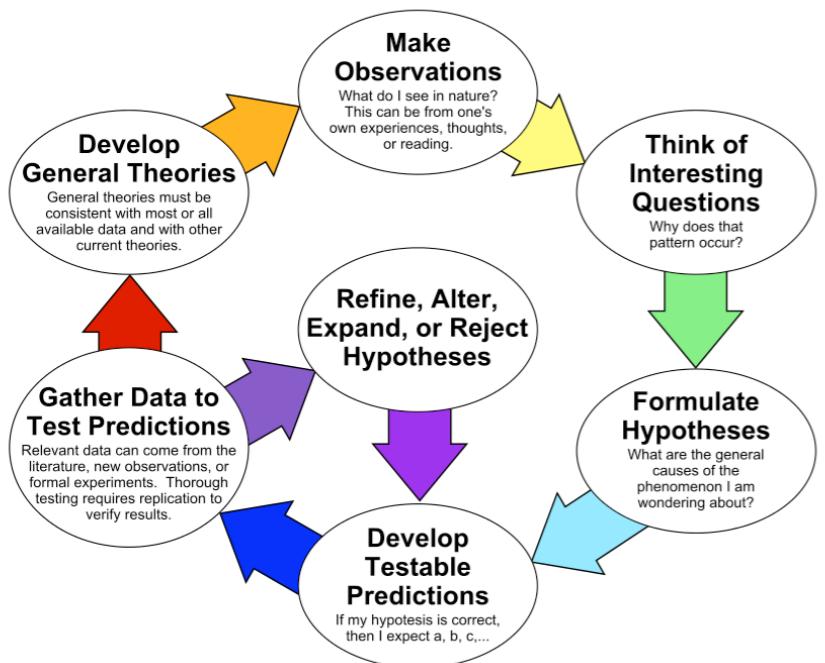
EDM  
Academic Analytics



Drew Conway  
“The Data Science  
Venn Diagram”  
(2010)

Educational Statistics  
Psychometrics  
Psychology

# Philosophy of Science Spectrum



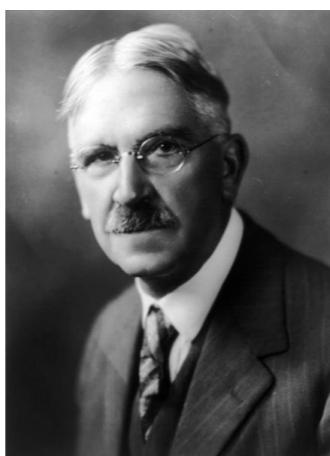
EDCT  
GE-2550



Mr Vargas



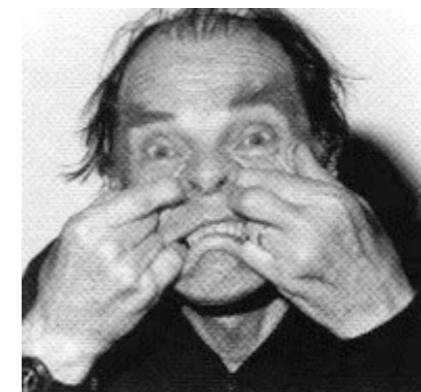
Delores  
Etter



John  
Dewey



William  
James



Paul  
Feyerabend



Morty

# Today

Introduction to the Field(s)

Candy

Course Details

Q&A

Quiz

Defining Educational Goals

Educational Goals

Analytic Strategy + Consent

Pros & Cons of variables

# Course Details

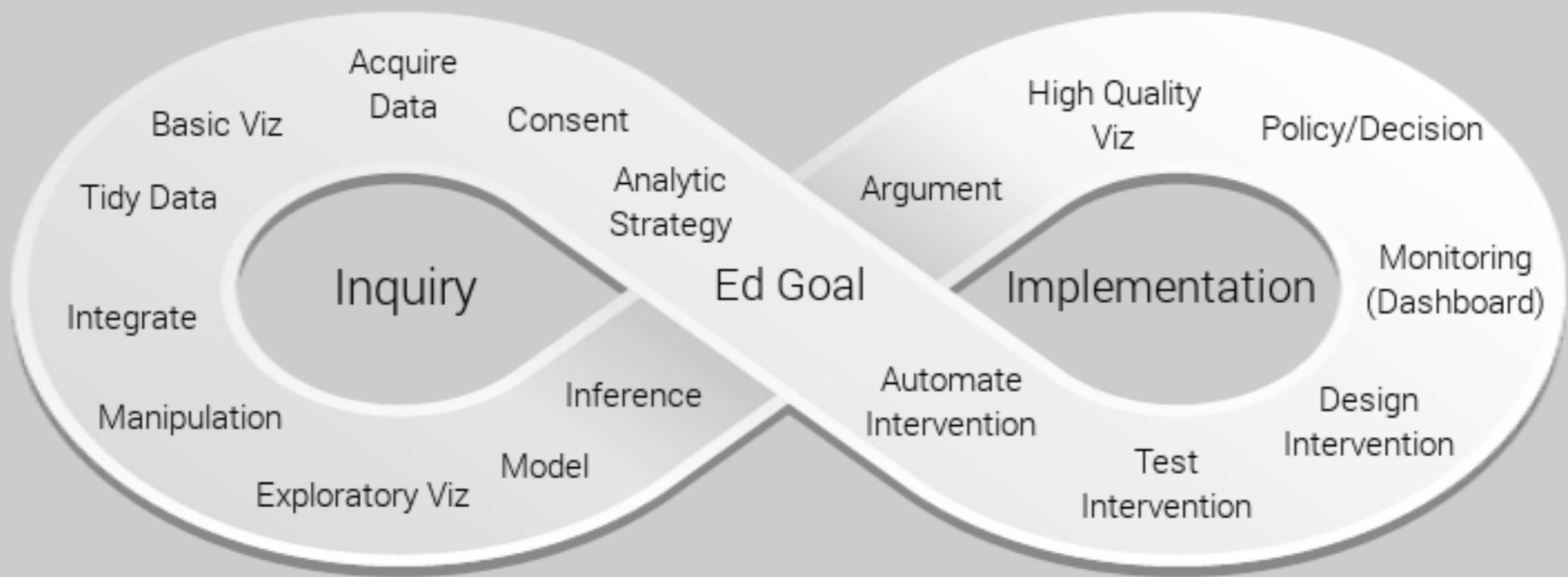
# Overall goal

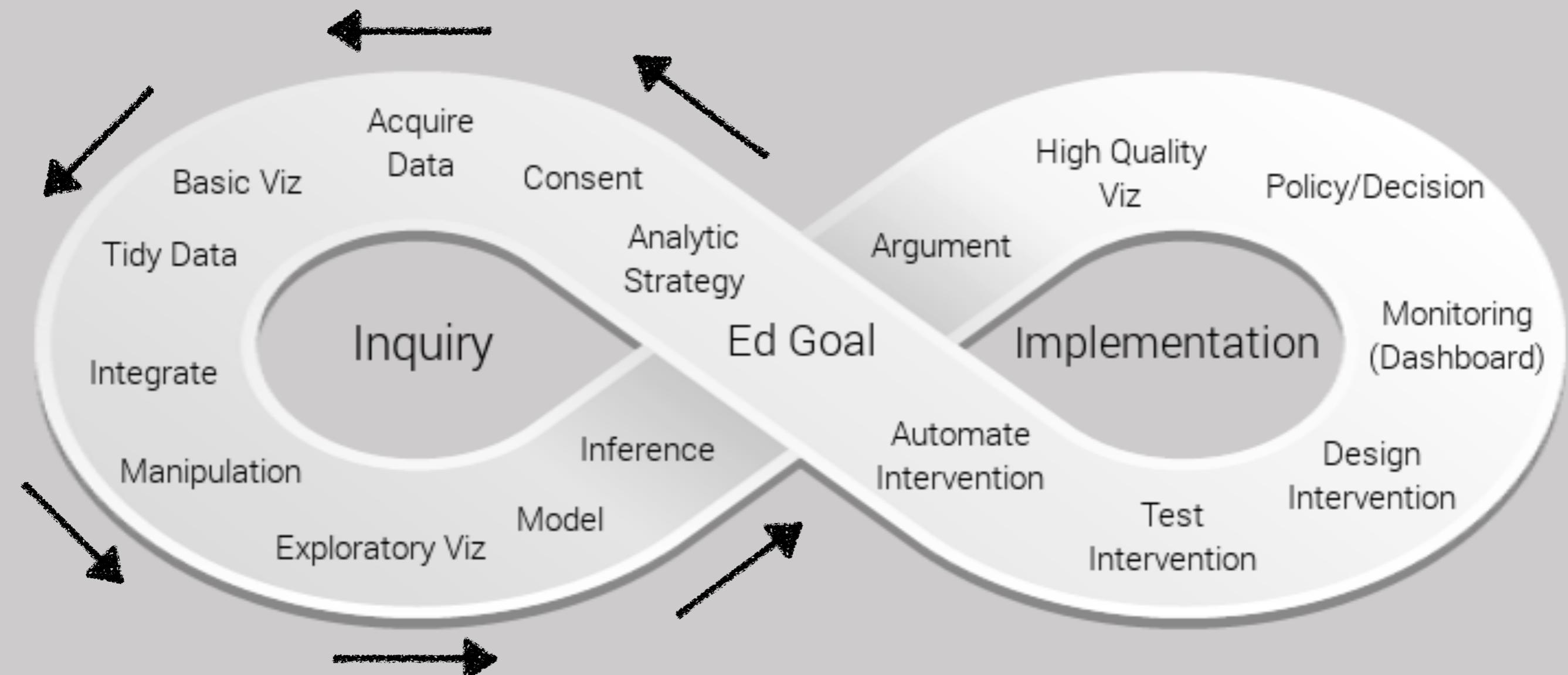
To investigate what the “Data Age” means for learning.

# Course Goals

| Short Term                                      |                                      | Long Term   |                              |
|---|--------------------------------------|---|------------------------------|
| <b><i>Content</i></b>                           | <b><i>Skills</i></b>                 | <b><i>Abstractions</i></b>  | <b><i>Habits</i></b>         |
| remember...,<br>understand...                   | demonstrate...                       | synthesize...,<br>argue...  | organize...,<br>implement... |
| Conceptual<br>basis of<br>methods, use<br>cases | R, git,<br>application of<br>methods | Evaluate<br>broader<br>implications,<br>have opinions,<br>methods<br>schema | Workflow,<br>documentation   |

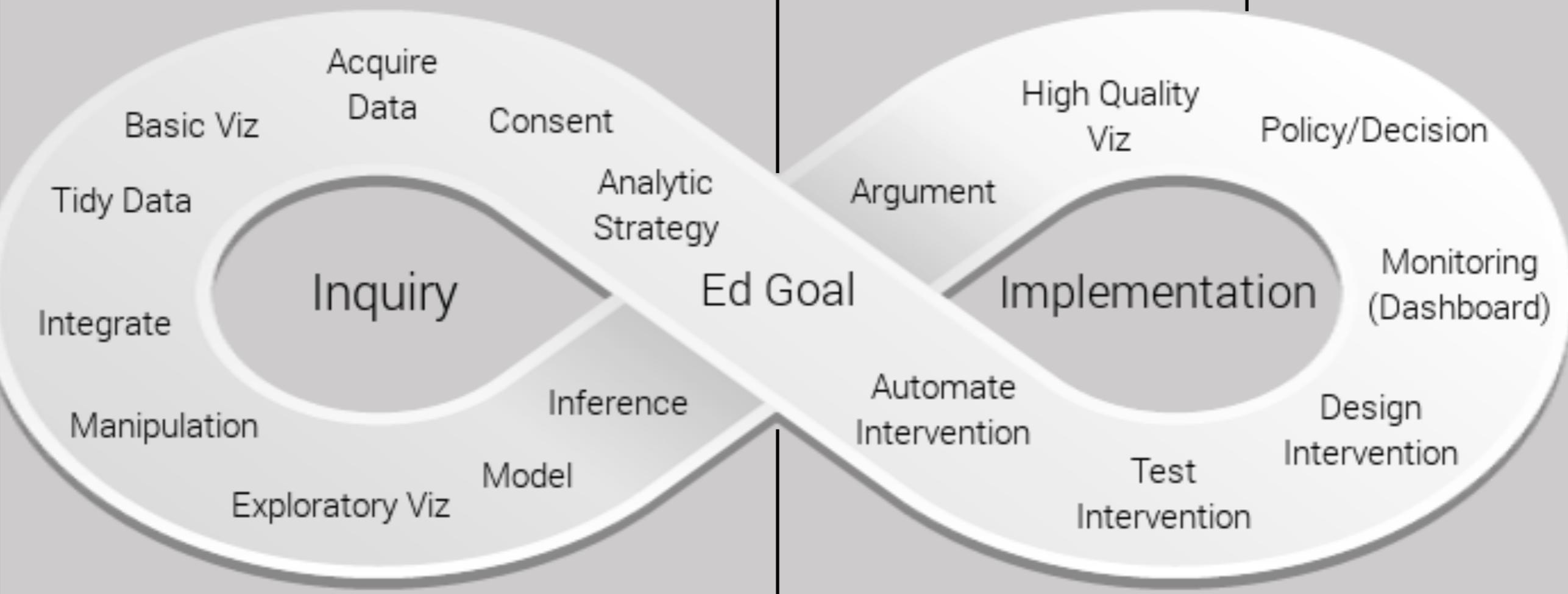
# Ed Data Science Cycle





Weeks 1-3

Weeks 12-14



Weeks 4-11

# Methods

- Data tidying
- Dimensionality Reduction
- Social Network Analysis
- Prediction
- Natural Language Processing

# What is the work?

- Attend class
- Weekly readings
- Comment on readings
- Weekly in class questionnaire
- Maintain documentation of work (Github, R Markdown, Zotero)
- Ask one question on Stack Overflow
- In person meeting with instructor
- 8 short assignments (including one group assignment)
- Group presentation of Assignment 8
- Produce one argument about learning using data from the class

# Assessment

***How to assess a course that is all about how difficult it is to measure learning?***

- Assess your preparedness to do the work after the course has ended
- Two measures:
  - Contribution: assignments, comments, quizzes
  - Organization: keep a record of what you have done (Zotero, Git, Markdown)

# Tools

- Git/Github



- R/RStudio



- Twitter



- Zotero



- Tools that are worth learning in & of themselves
- Tools that we can easily extract data from
- Tools that you can use unrestricted in the future

# Git/Github



- Git is a common version control system
- Github is an online hosting service for Git files

We will be using Git/Github to:

- Keep track of work
- Document problems
- As an LMS (assignment submission, etc)



- R is a scientific programming language
- RStudio is an Integrated Development Environment

# Why use R?

## Pros

- Free
- Platform independent
- Actively developed by a large community of users
- Functionality is VAST
- Help resources are VAST
- Best graphics (At the moment)
- All the cool kids use it

# Why use R?

## Cons

- Slow (for some tasks)
- Not designed to build tools
- Relies on vectorization
- Non-intuitive structures & syntax with respect to other programming languages

# Why use R?

It is designed for people who are *learning* a method as *they apply it*.

# Twitter (#DSE16)



Three main uses for Twitter:

- Crowd source answers to your questions
- Comment on readings
- Backchannel in class (#IDKWGO)

# Zotero



- Bibliographic software
- Lives in browser (where you do research)
- Open source

# Data

As much as possible we will be generating the data we analyze:

- Gain the experience of being on both sides of the data collection & analysis
- Easier to be invested in the content

# But...

- Need to separate what we consider “assessment” and what we consider “activity”
- Need to put limitations on use (consent, time limit)
- Need to develop trust

# Last word

If you are having trouble doing the work for whatever reason do not hesitate to reach out (the sooner the better):

[charles.lang@nyu.edu](mailto:charles.lang@nyu.edu)

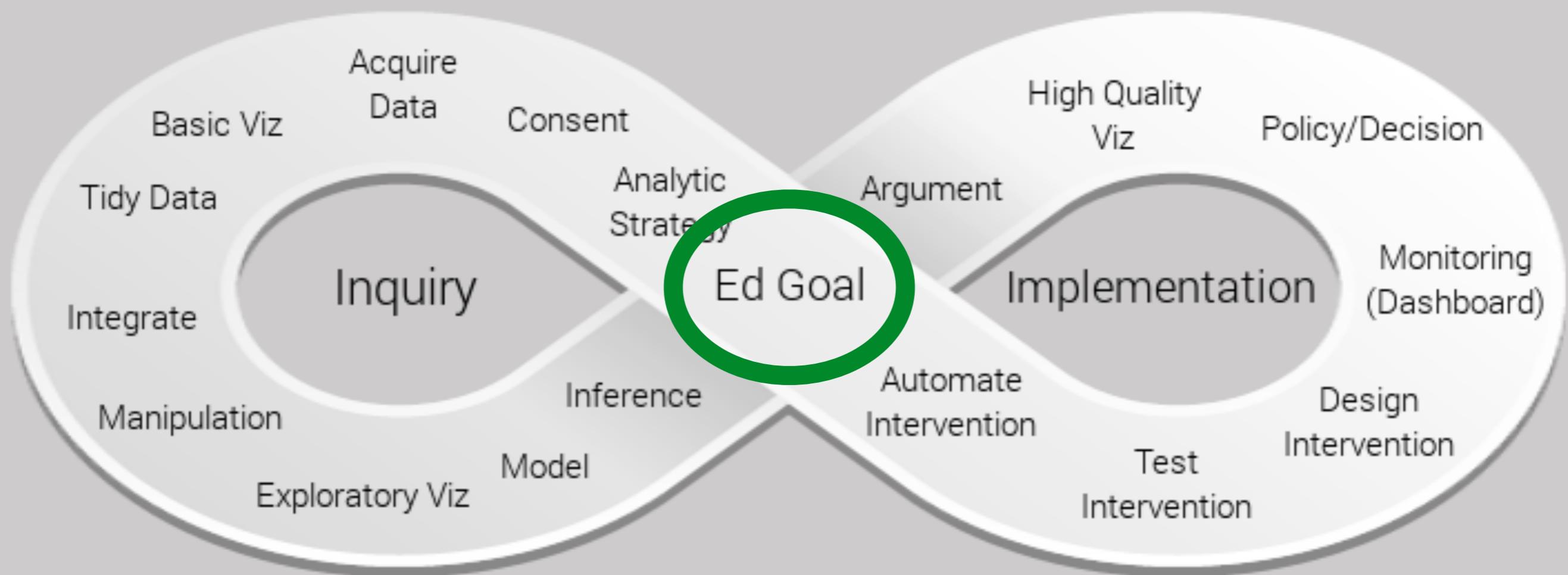
@learng00d

Syllabus

<https://github.com/data-science-in-ed/Syllabus>

# Educational Goals

# Ed Data Science Cycle



“It’s better to be at the bottom of the ladder you want to climb than at the top of the one you don’t.”

–Stephen Kellogg, 2013

# Google Books Ngram Viewer



# Why?

- Education is controversial
- Concepts cross disciplinary lines
- Vocabulary is slippery
- Systems are complex

# Characteristics

- Specific/tightly defined
  - Defined level of inference
- No jargon
  - Contextual framework
- Measurable
  - They can change (BUT keep a record of the change)
- Attainable
  - Can be to answer a question (To determine the time of day students write the most words.)
- Relevant
  - Needs to be workshopped with others
- Time-limited

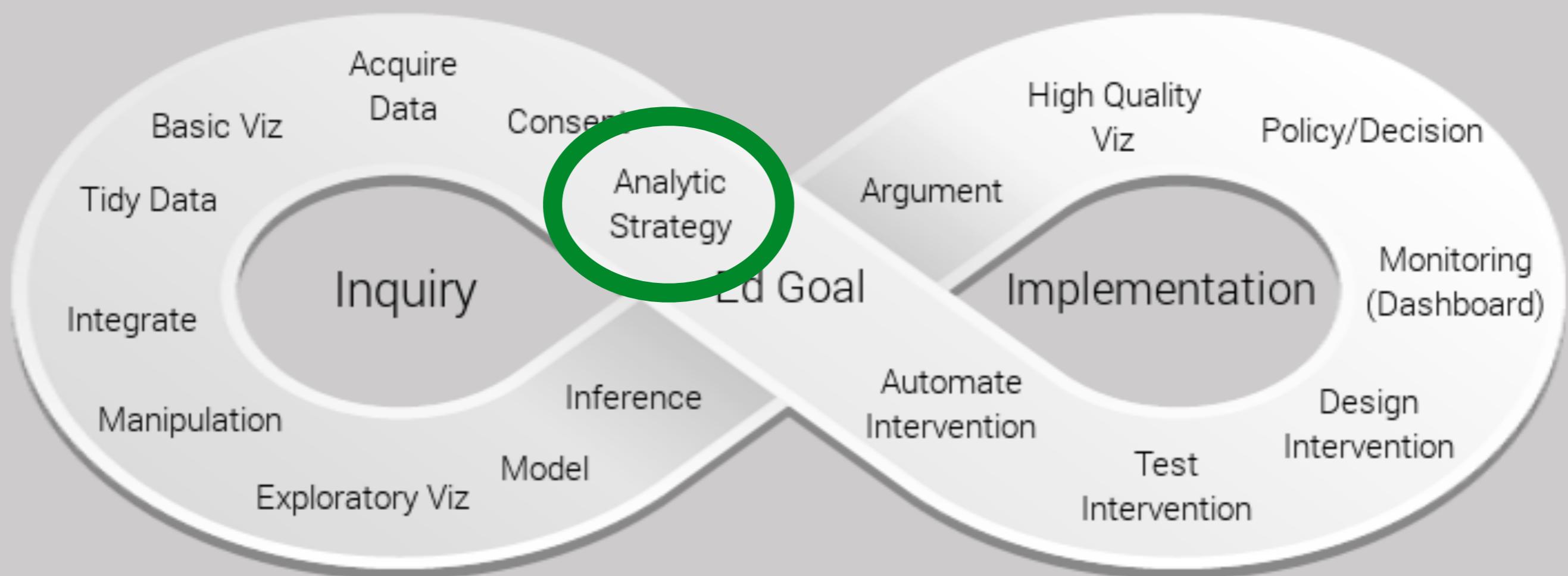
# Activity 1

1. Find a partner
2. Determine one educational goal that you would like to define for yourself, for this semester
3. What information would you need to allow you to determine if you had met that goal

# Analytic Strategy

1/30/16 4:39 PM

# Ed Data Science Cycle

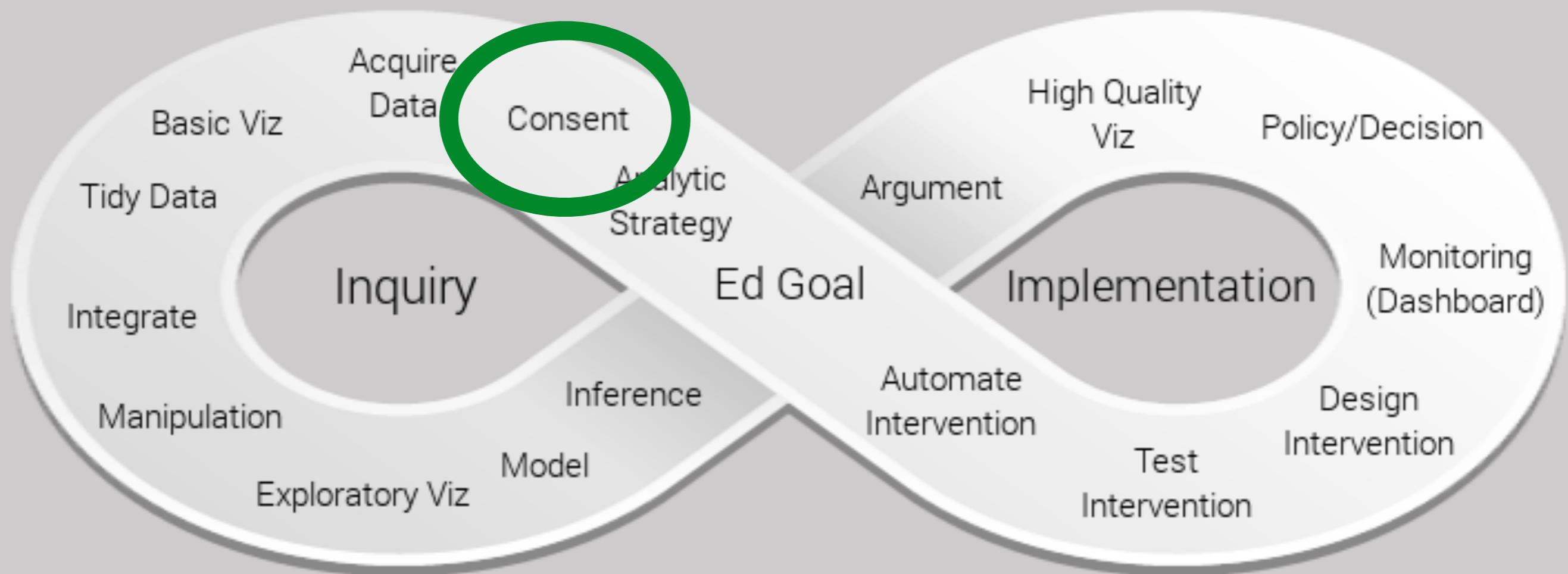


# Analytic Strategy

- Going to use a range of methods
- Data will be erased at the end of the course



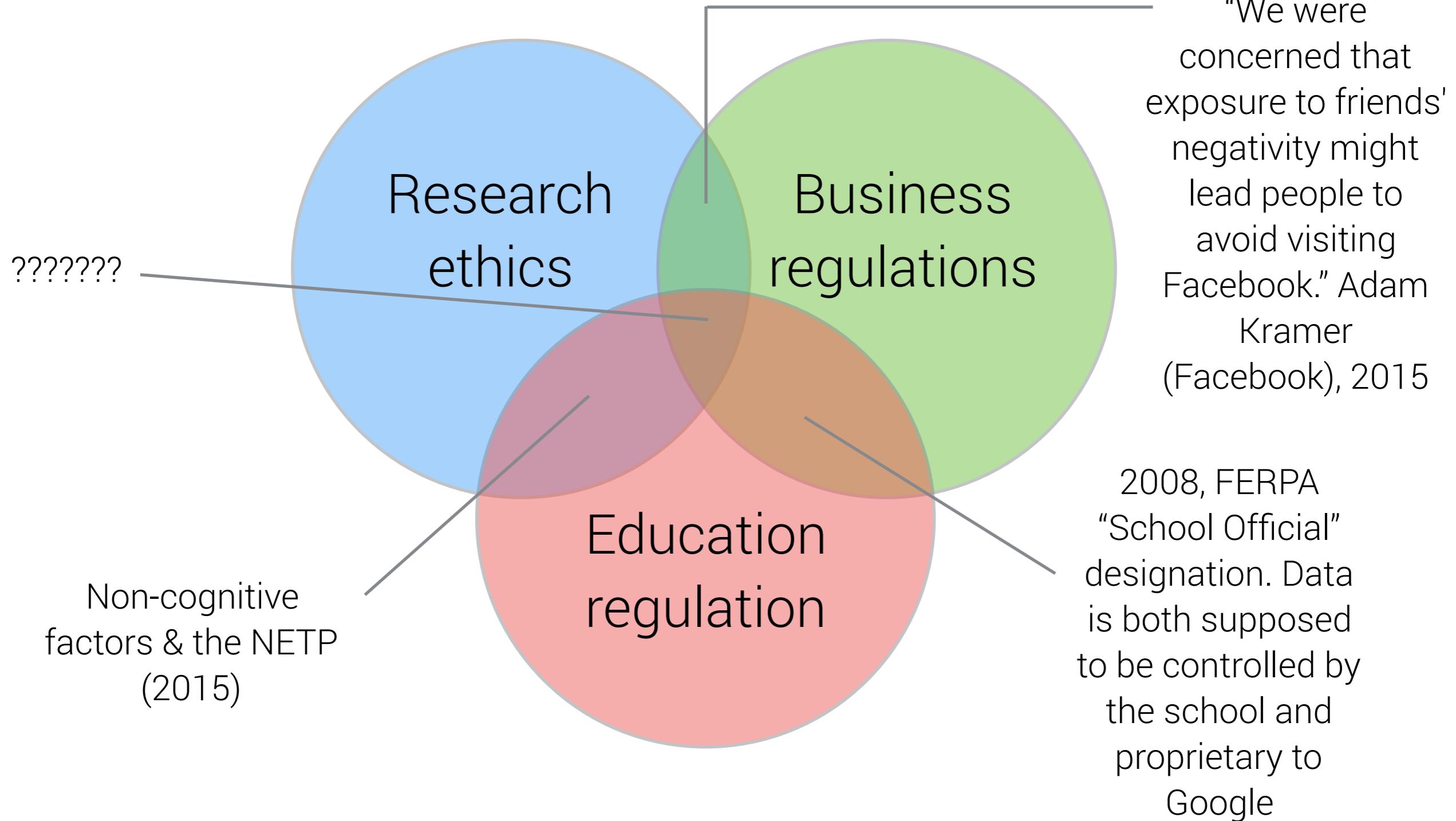
# Ed Data Science Cycle



# Consent

1. Disclosure
2. Capacity
3. Voluntariness

# Class of Civilizations



What are people used  
to?

# Cognitive Dissonance

- Transparency vs. obfuscation
- Siloed information vs. linked functionality
- Defined use vs. mine for insight



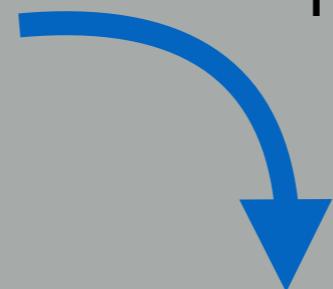
# Transparency vs. Obfuscation

| Terms of Service   | Privacy Policy | Community Guidelines |
|--|----------------|----------------------|
| <p>No individual under the age of thirteen (13) may use the Services, provide any personal information to Tumblr, or otherwise submit personal information through the Services (including, for example, a name, address, telephone number, or email address). You may only use the Services if you can form a binding contract with Tumblr and are not legally prohibited from using the Services.</p> <p>You have to be at least 13 years old to use Tumblr. We're serious: It's a hard rule, based on U.S. federal and state legislation. "But I'm, like, 12.9 years old!" you plead. Nope, sorry. If you're younger than 13, don't use Tumblr. Ask your parents for a Playstation 4, or try books.</p> |                |                      |

# Siloed Information vs. Linked Functionality



A screenshot of the Syllabus Finder website. At the top is a black and white photo of a person in a hoodie waving. Below it is the title "Syllabus Finder". A search bar contains the text "discala spencer" and a "search" button. Below the search bar is the text: "Searching 1,121,847 syllabi at the Center for History and New Media and over 500,000 syllabi via Google". At the bottom is a link: "From Dan Cohen's blog: the most popular syllabi in [history](#) and [philosophy](#)".



Generation of  
“Popularity” Indexes

(2011)



A screenshot of the Open Syllabus Explorer website. At the top right are links for "About", "People", and social media icons for Twitter, Facebook, and GitHub. In the center is the title "Open Syllabus Explorer" in large white letters, with a "beta 0.4" badge to its right. Below the title is the tagline: "Mapping the college curriculum across 1M+ syllabi." The background is a blurred image of a classroom.

(2016)

# Defined Use vs. Mine for Insight

- Restrict use = lose insight
- “Educational use” might be useless
- EG - Time limits on user profiles

# Activity 2

1. Get into candy-based groups (Snickers, Milky Ways, etc)
2. Each member of the group is assigned a role:
  - Student = snickers, teacher = MW midnight, parent = MW Regular, ed tech CEO = 3 Musk, researcher = Twix
3. For your given role discuss pros & cons for students allowing the type of data assigned to your group to be collected in an educational setting for the purpose of *improving learning*

# Activity 2

- 1: Environment measures (eg. room temperature)
- 2: Formative assessments (eg. weekly quizzes)
- 3: Location data (eg. location in building)
- 4: Biometric data (eg. Fitbit data)
- 5: External social website data (eg. Twitter/Facebook)
- 6: Demographic information (eg. gender, age)