

What is Data Science?

January 24, 2019

Data Science CSCI 1951A

Brown University

Instructor: Ellie Pavlick

HTAs: Wennie Zhang, Maulik Dang, Gurnaaz Kaur

Waitlist

- If you are not registered, make sure you are on the waitlist (link on CAB)
- We have a *little* wiggle room in the enrollment cap
- Indicate relevant extenuating circumstances, we will try to prioritize fairly

What is Data Science?

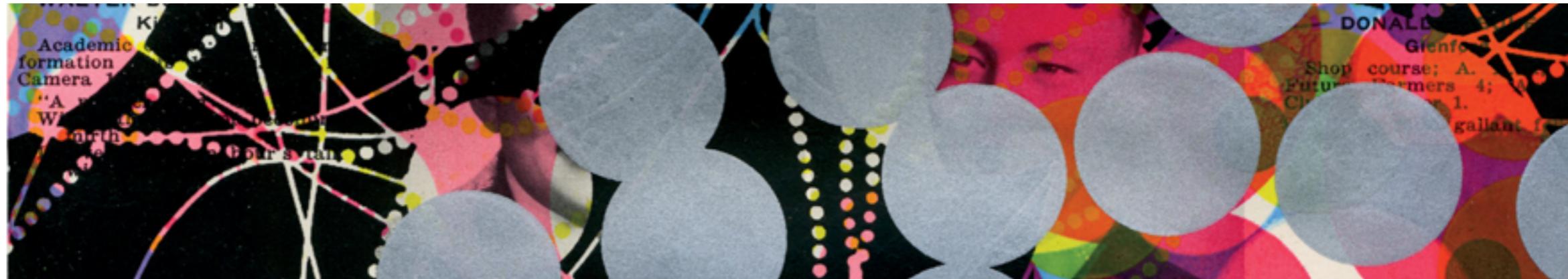


DATA

Data Scientist: The Sexiest Job of the 21st Century

More than anything, what data scientists do is make discoveries while swimming in data. It's their preferred method of navigating the world around them. At ease in the digital realm, they are able to bring structure to large quantities of formless data and make analysis possible. They identify rich data sources, join them with other, potentially incomplete data sources, and clean the resulting set. In a competitive landscape where challenges keep changing and data never stop flowing, data scientists help decision makers shift from ad hoc analysis to an ongoing conversation with data.

Data scientists realize that they face technical limitations, but they don't allow that to bog down their search for novel solutions. As they make discoveries, they communicate what they've learned and suggest its implications for new business directions. Often they are creative in displaying information visually and making the patterns they find clear and compelling. They advise executives and product managers on the implications of the data for products, processes, and decisions.



DATA

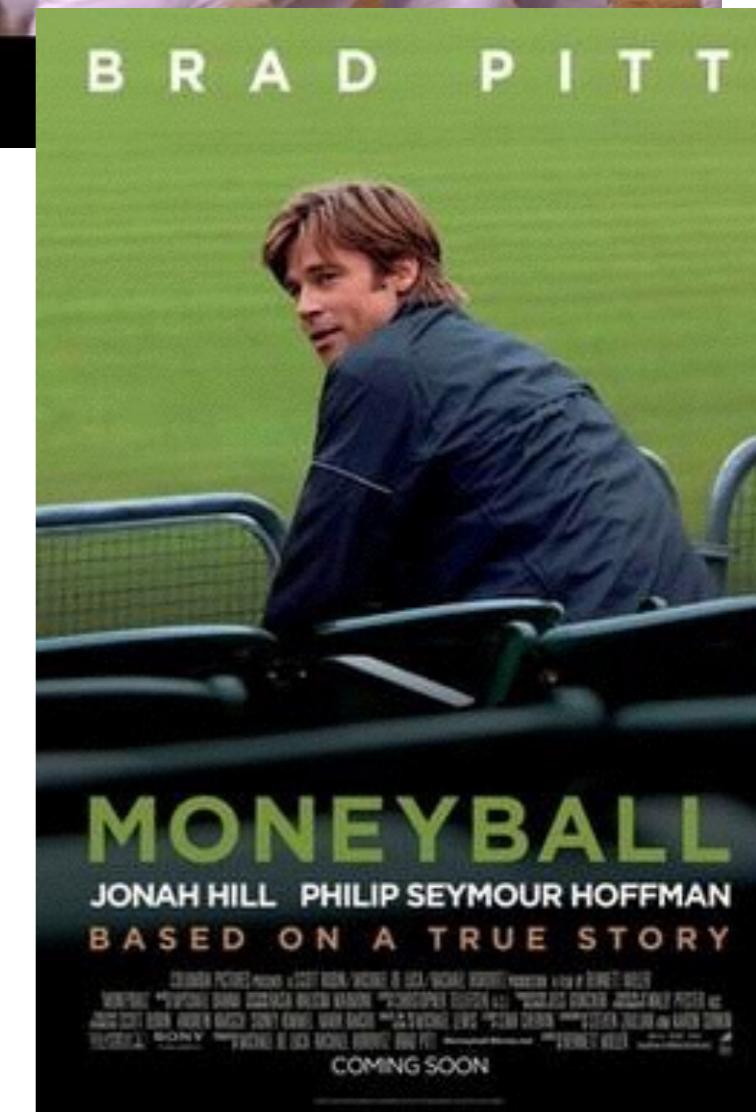
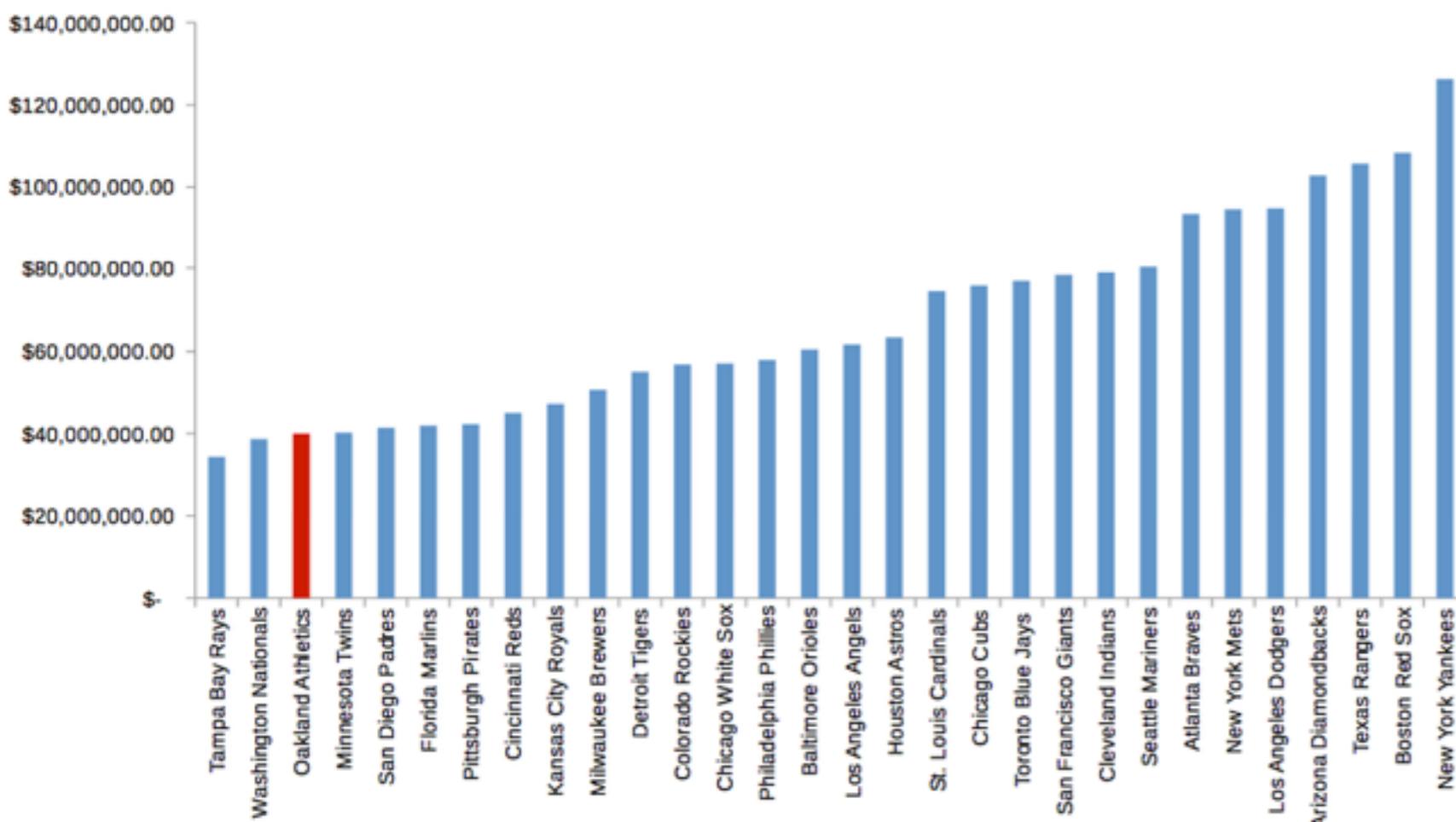
Data Scientist: The Sexiest Job of the 21st Century

More than anything, what data scientists do is make discoveries while swimming in data. It's their preferred method of navigating the world around them. At ease in the digital realm, they are able to bring structure to large quantities of formless data and make analysis possible. They identify rich data sources, join them with other, potentially incomplete data sources, and clean the resulting set. In a competitive landscape where challenges keep changing and data never stop flowing, data scientists help decision makers shift from ad hoc analysis to an ongoing conversation with data.

Data scientists realize that they face technical limitations, but they don't allow that to bog down their search for novel solutions. As they make discoveries, they communicate what they've learned and suggest its implications for new business directions. Often they are creative in displaying information visually and making the patterns they find clear and compelling. They advise executives and product managers on the implications of the data for products, processes, and decisions.

Moneyball!

Moneyball Year (2002)
MLB Team Salaries



Obama Campaign



CONTROL

A screenshot of a campaign website for Barack Obama and Joe Biden. The top navigation bar features the "OBAMA BIDEN" logo. Below it, a banner reads "DINNER WITH BARACK" with the subtitle "Your chance to meet the President". A red "GET STARTED" button is visible. The main image shows a smiling Barack Obama in a white shirt and red tie, standing next to a woman whose back is to the camera. The background is a dark room with a piano. The text "DINNER WITH BARACK" is overlaid on the left side of the image. Below the image, the text "YOU'RE INVITED. WE'LL COVER YOUR AIRFARE." is displayed. At the bottom of the page, there is a section titled "About this event" with a "Read more" link, followed by a "DONATE" button and the "OBAMA BIDEN" logo.

IMAGE VARIATION

A screenshot of the official website for the "Dinner with Barack" event. The top navigation bar features the "OBAMA BIDEN" logo. Below it, the main title "DINNER WITH BARACK" is displayed in large, bold letters, with the subtitle "Your chance to meet the President" underneath. A red "GET TICKETS" button is positioned to the right. The central image shows a group of diverse individuals seated around a long table, engaged in conversation. To the right of the photo, the text "DINNER WITH BARACK" is repeated in a decorative font, followed by "You're invited." and "How it covers your airfare." Below the photo, there is descriptive text about the event's purpose and how it covers travel expenses. At the bottom of the page, the "OBAMA BIDEN" logo is centered, along with links for "About Us," "Privacy Policy," and "Terms of Service." A note about contributions and gift tax information is also present.

↑ +19%

Just as the beam is

CONTROL


Logout | Create account



You could be there

Join our team. Barack Obama, Sheila Jackson, Eric Holder, Anthony, and many, many more influential people are all cheering us on to support President Obama.

You are a great role model for the President – and should bring with you one of the best basketball players ever.

[More ways to be involved](#)

Contributor

Select amount

Credit card



Make this a monthly donation to sustain this campaign. In
the future see more like this. (Check this box to receive more info)

Employment

Postal law requires us to use our best efforts to collect and report the name, mailing address, occupation, and employer of individuals whose contributions exceed \$200 in an election cycle.



"SEQUENTIAL"

A screenshot of a campaign donation page for Obama Biden. The top navigation bar features the "OBAMA BIDEN" logo and links for "Sign In" and "Create account". Below the header is a blue banner with the "OBAMA CLASSIC" logo and a call to action: "You could be there" followed by a list of names: "Barack Obama, Mattie Edging, Diane Sawyer, Common, Audra McDonald, and many more...". It also mentions "This will be a great way to raise money for President Obama—and spend time with some of the most interesting people still..." and "More details on the announcement". The main content area has a dark blue background with white text asking "How much would you like to donate today?". It includes a "Select amount" dropdown with options: "\$15", "\$25", "\$50", "\$100", "\$250", "\$500", "\$1,000", and "Other amount". A large "CONTINUE" button is at the bottom. At the very bottom, there's a footer with links: "OBAMA BIDEN", "About", "Political Action Committee", and "Contact Us". On the right side of the page, there's a large image of Barack Obama shooting a basketball.

↑ +5%

Google's “40 Shades of Blue”



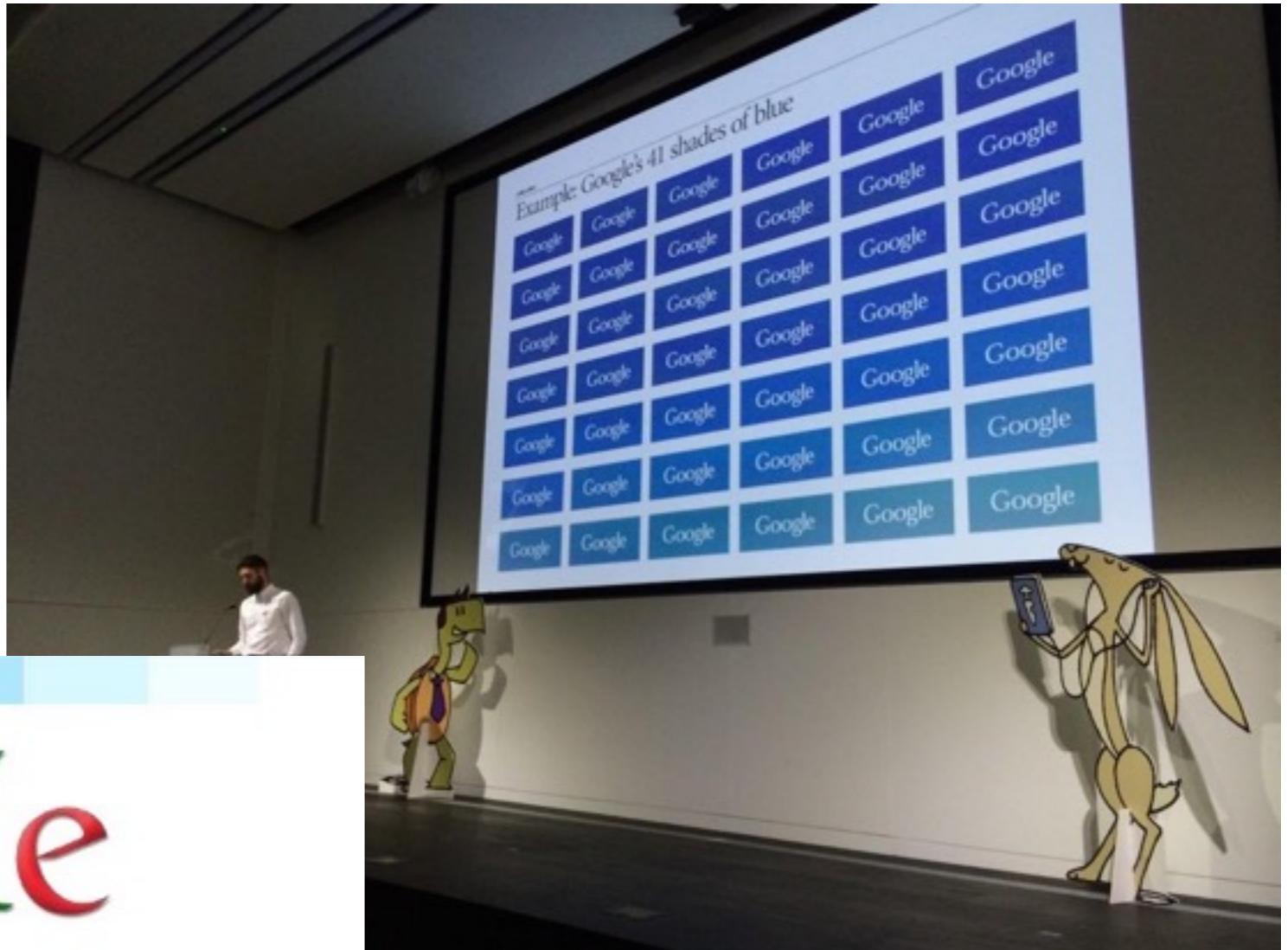
Google

a team at Google couldn't decide between two blues

they tested **41 shades** between each blue,
showing each one to 1% of their visitors
to see which one performs better

\$200 million of benefits

Why Google has 200m reasons to put engineers over designers. The Gaurdian.
The Origin of A/B Testing. Nicolai Kramer Jakobsen.





Data Science = Magic





LiveSlides web content

To view

Download the add-in.

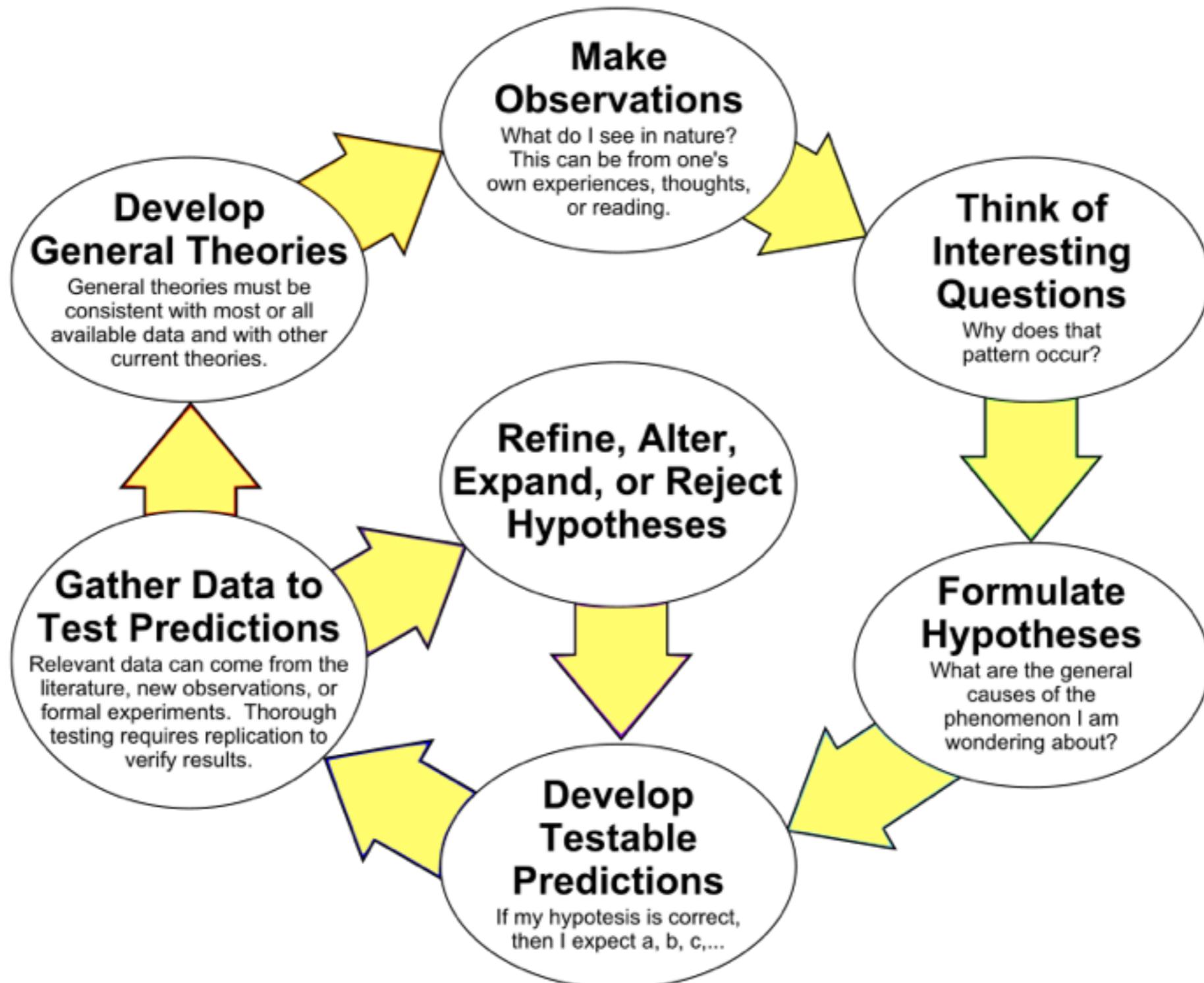
liveslides.com/download

Start the presentation.

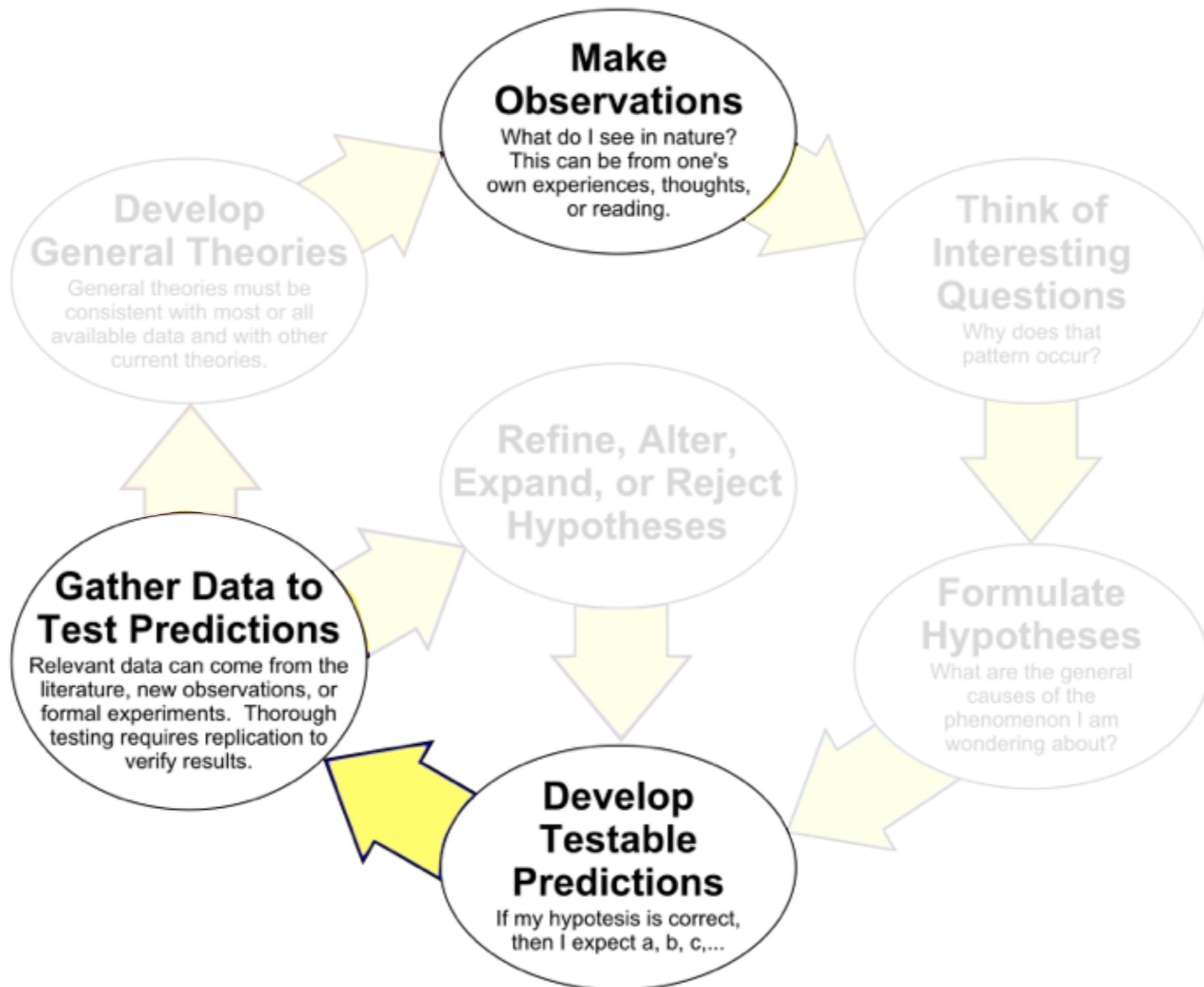
MACHINE LEARNING



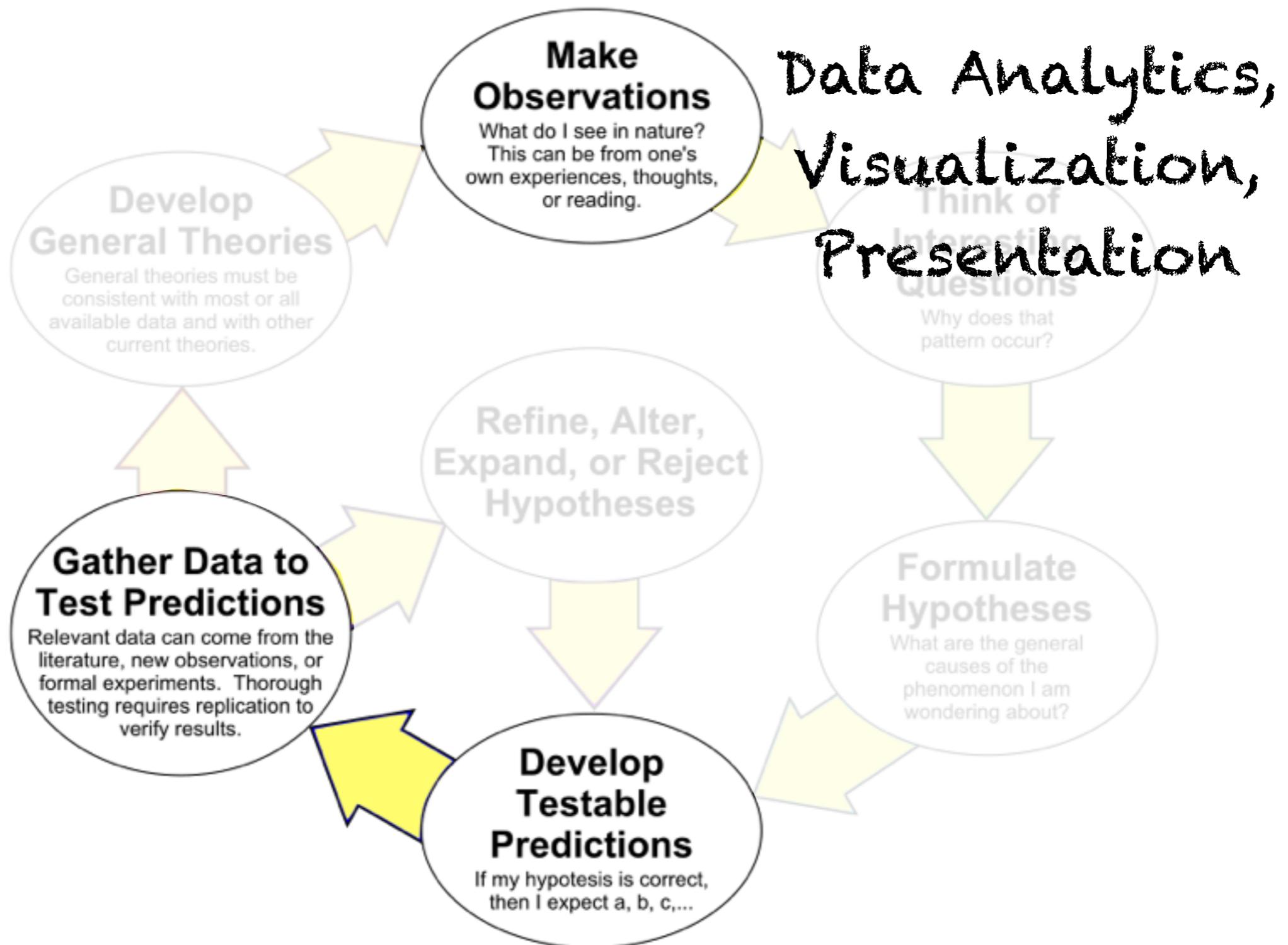
The Scientific Method



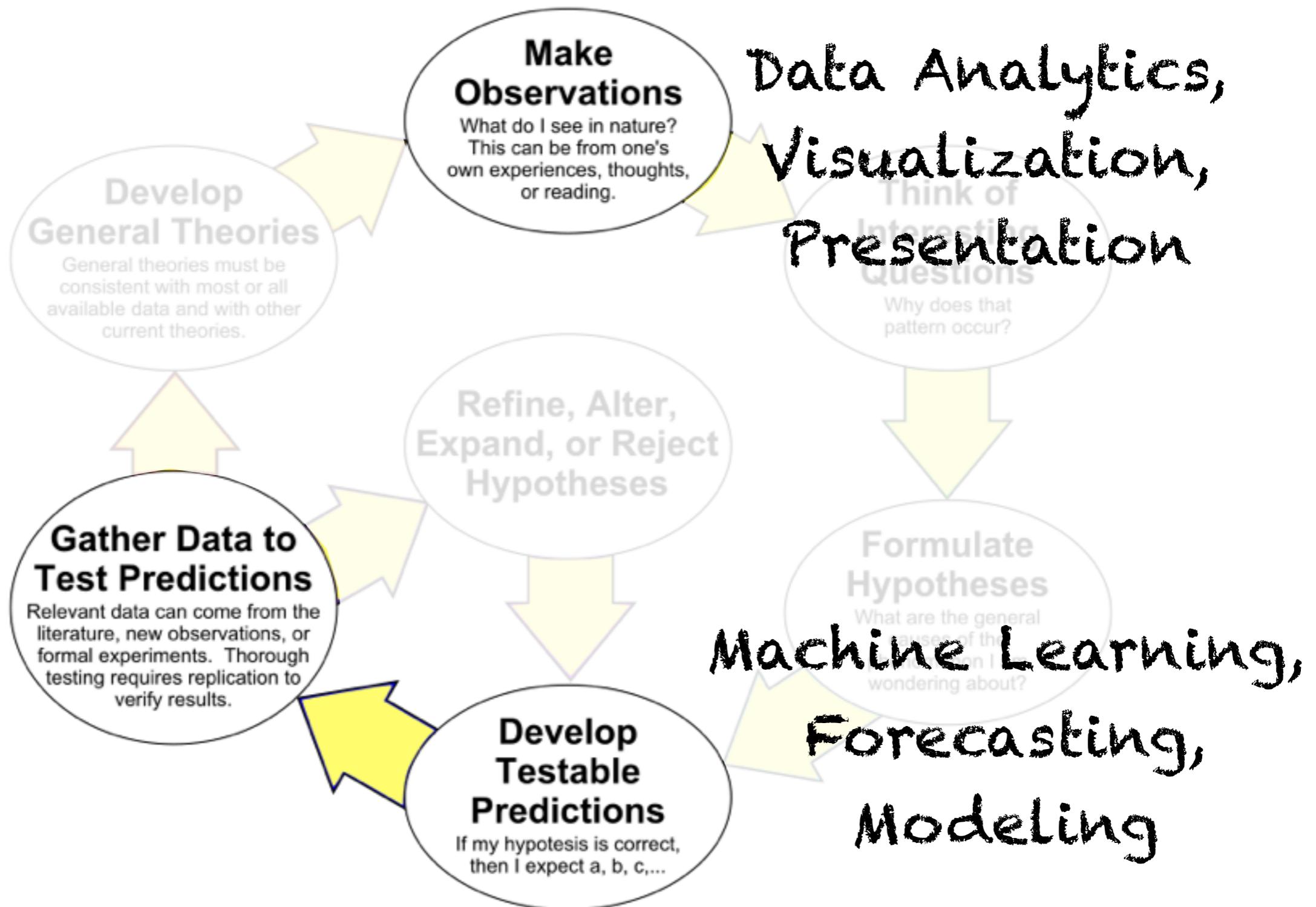
The Scientific Method



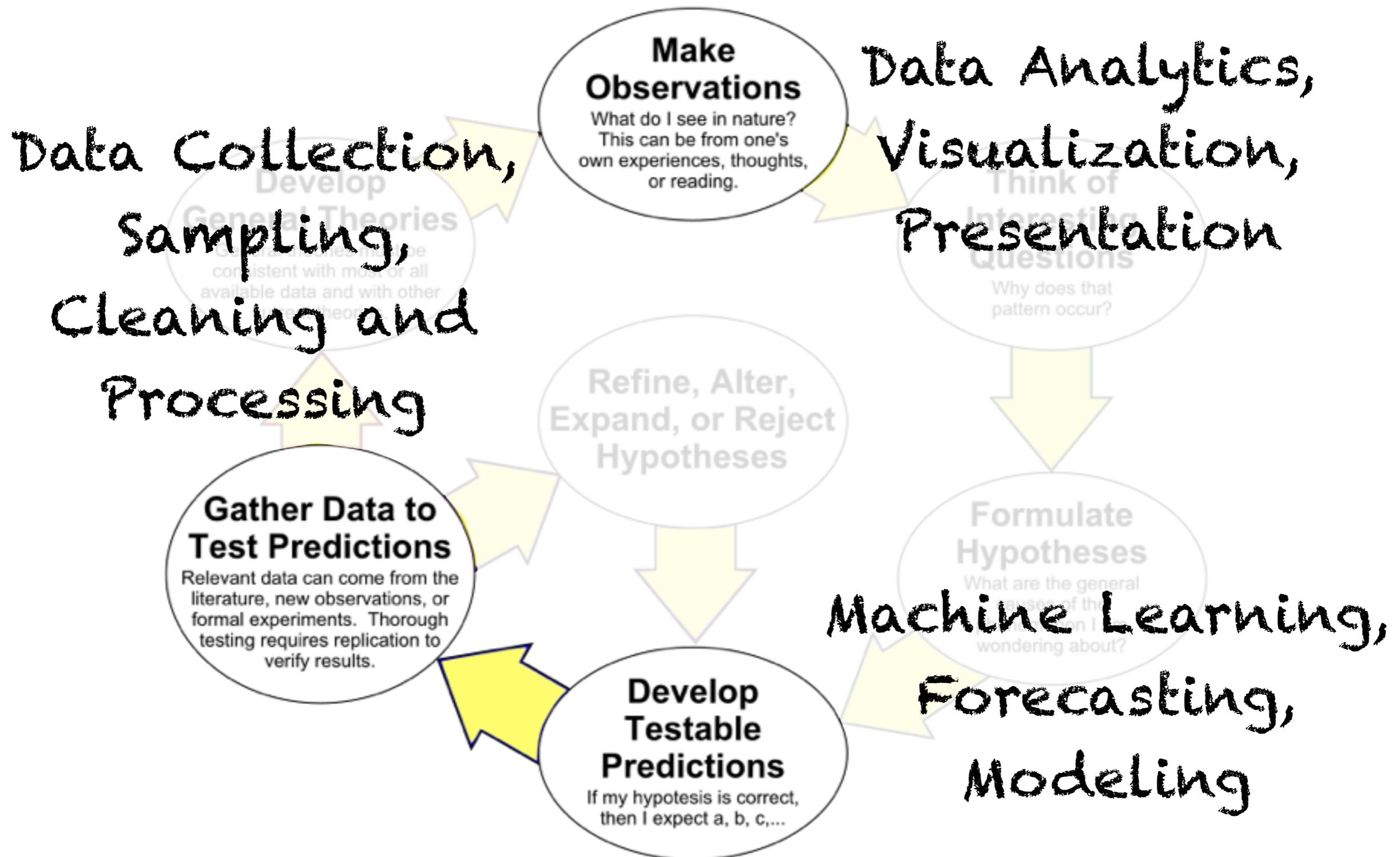
The Scientific Method



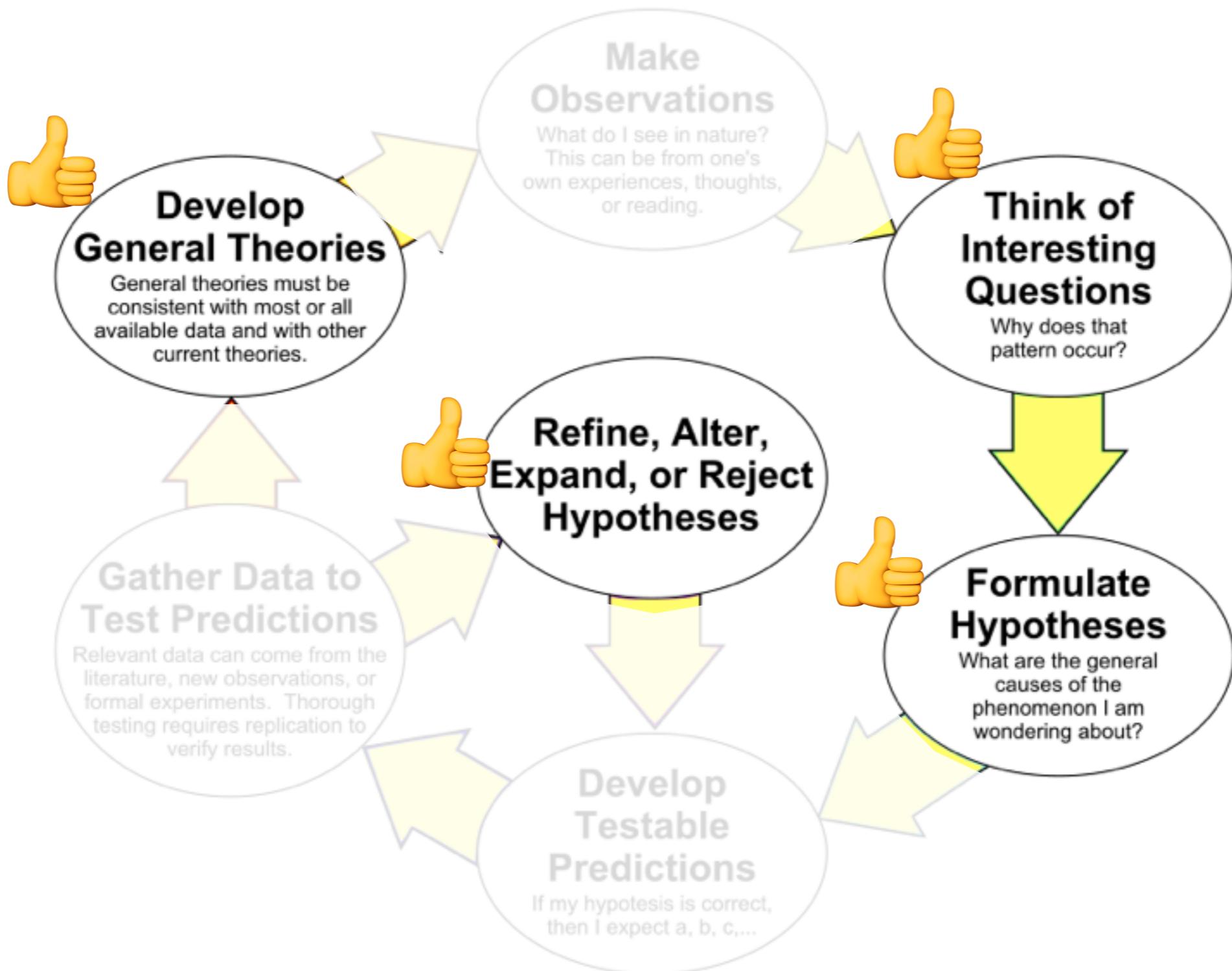
The Scientific Method



The Scientific Method



The Scientific Method



What is Data Science?





What is Dad's Dreaded Home Maintenance?



**Hang the
Christmas lights...**

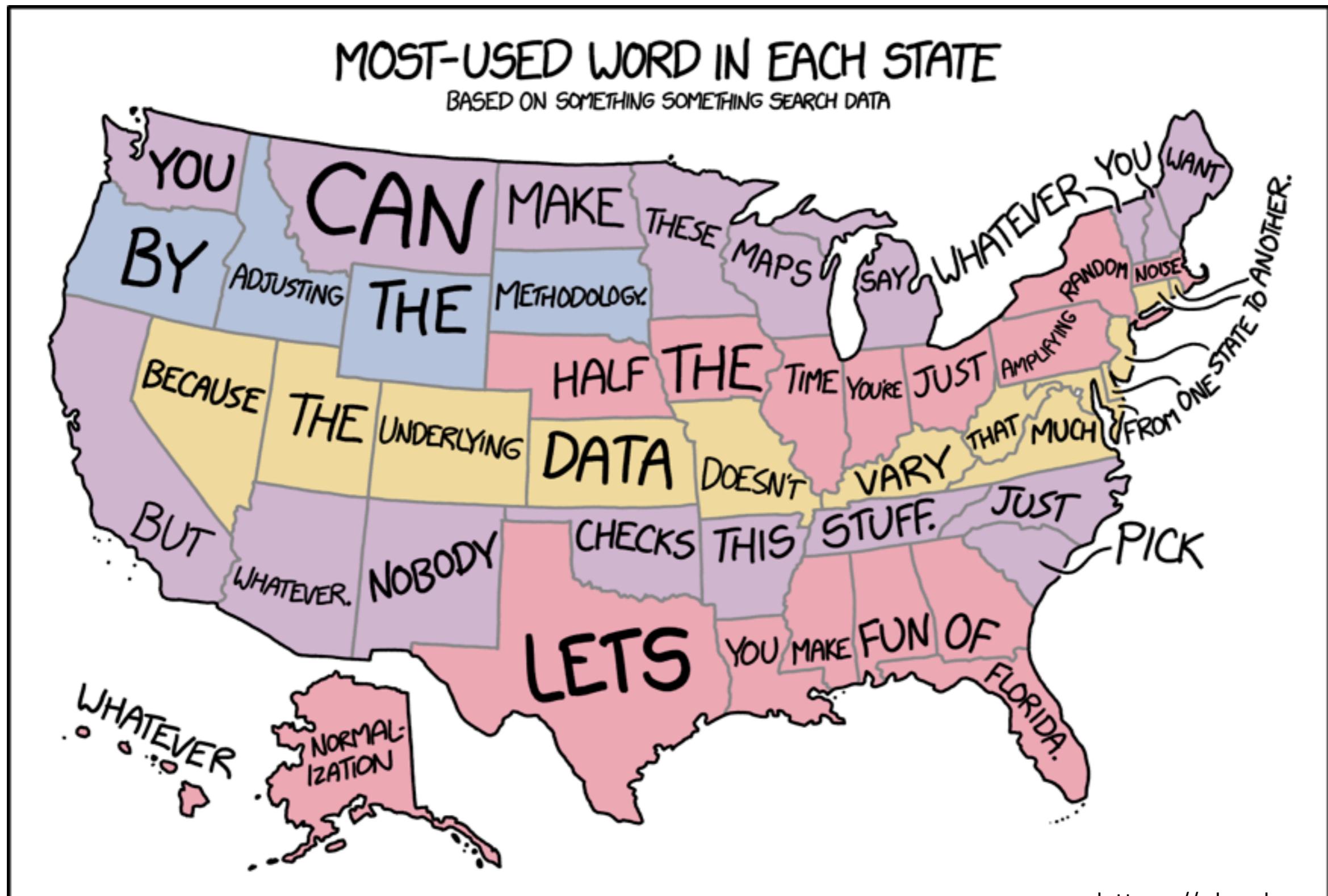
Data “Science”

Data “Science”



<https://www.dailydot.com/unclick/state-googled-2017>
<http://nerdgeeks.co/us-state-words-map>

Data “Science”



Data “Science”

- To be fair...

Data “Science”

- To be fair...
 - Intuition plays a huge role in the scientific method (“make observations” is Step 1).

Data “Science”

- To be fair...
 - Intuition plays a huge role in the scientific method (“make observations” is Step 1).
 - Exploratory analysis is necessary, its okay to not be all rigor all the time

Data “Science”

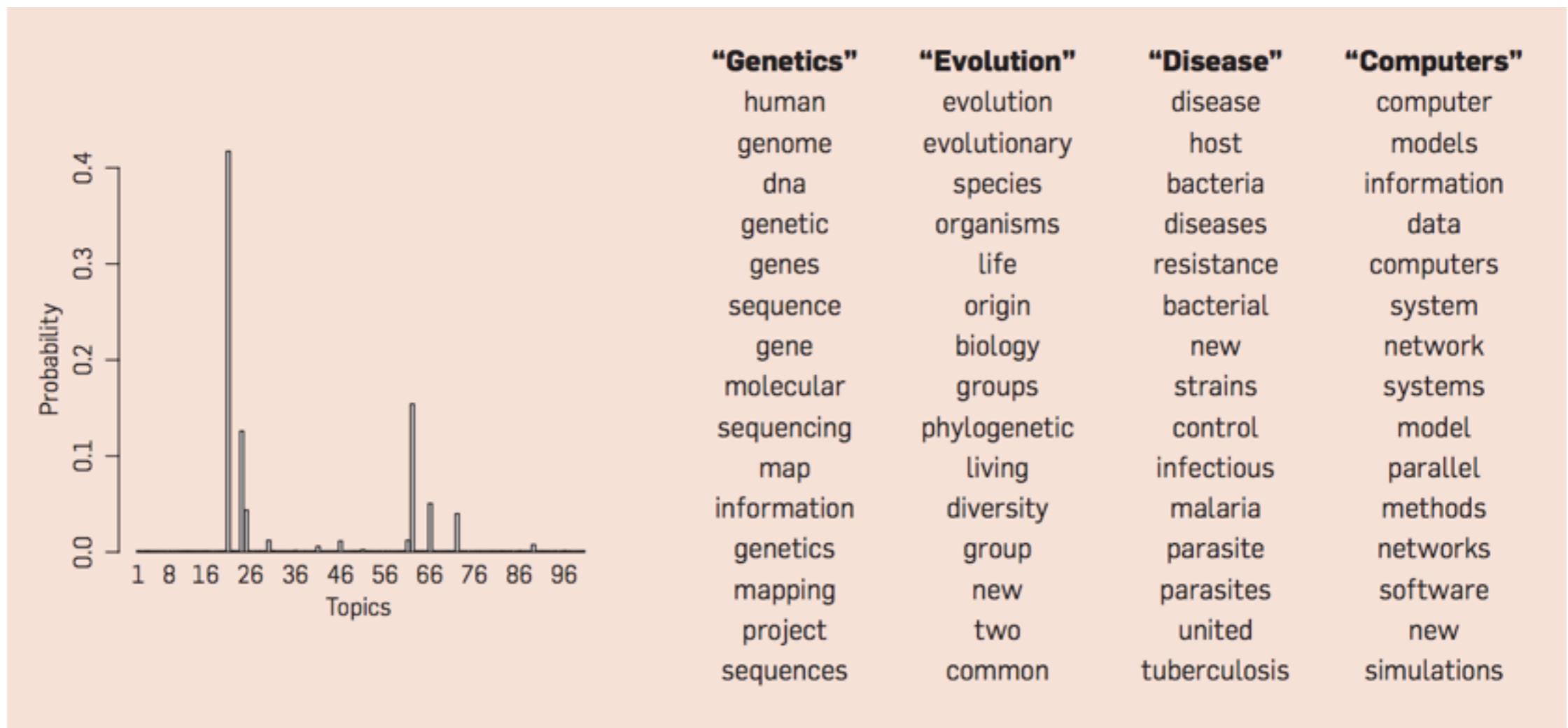
“Eyeballing it”



Facebook posts by age group

Data “Science”

“Eyeballing it”

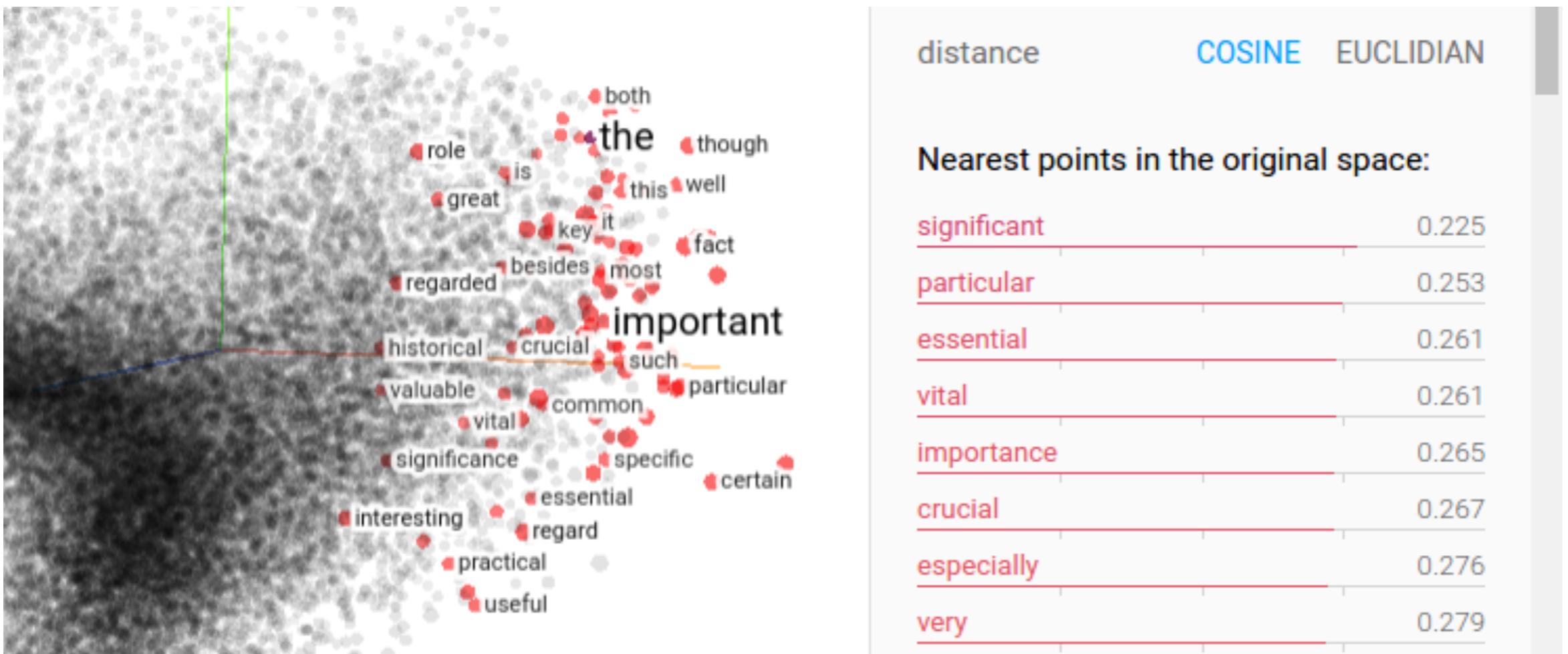


Frequent topics observed in 17,000 Science articles

Probabilistic Topic Models. Blei (2012).

Data “Science”

“Eyeballing it”



Similarity of words according on word2vec model

Data “Science”

- To be fair...
 - Intuition plays a huge role in the scientific method (“make observations” is Step 1).
 - Exploratory analysis is necessary, it’s okay to not be all rigor all the time
- But!

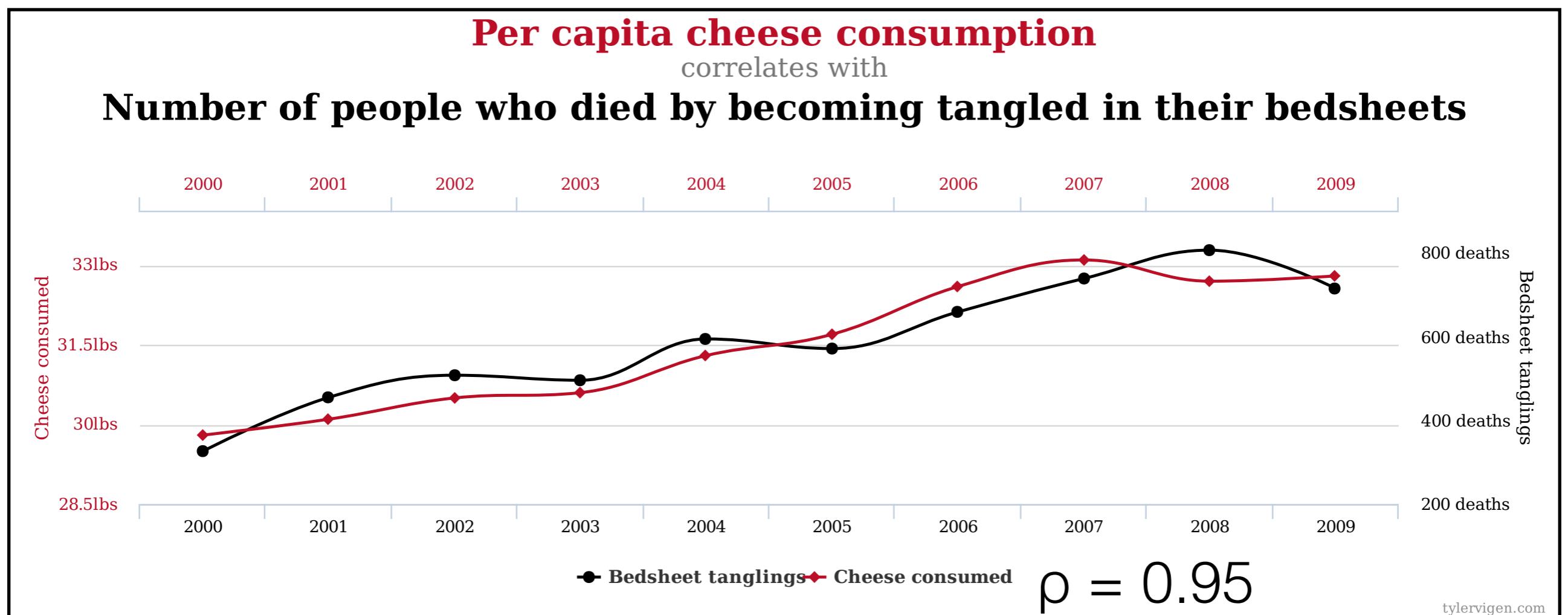
Data “Science”

- To be fair...
 - Intuition plays a huge role in the scientific method (“make observations” is Step 1).
 - Exploratory analysis is necessary, it's okay to not be all rigor all the time
- But!
 - Exploratory analysis (even when it involves the biggest of data) is meant to **form** a hypothesis, not test one

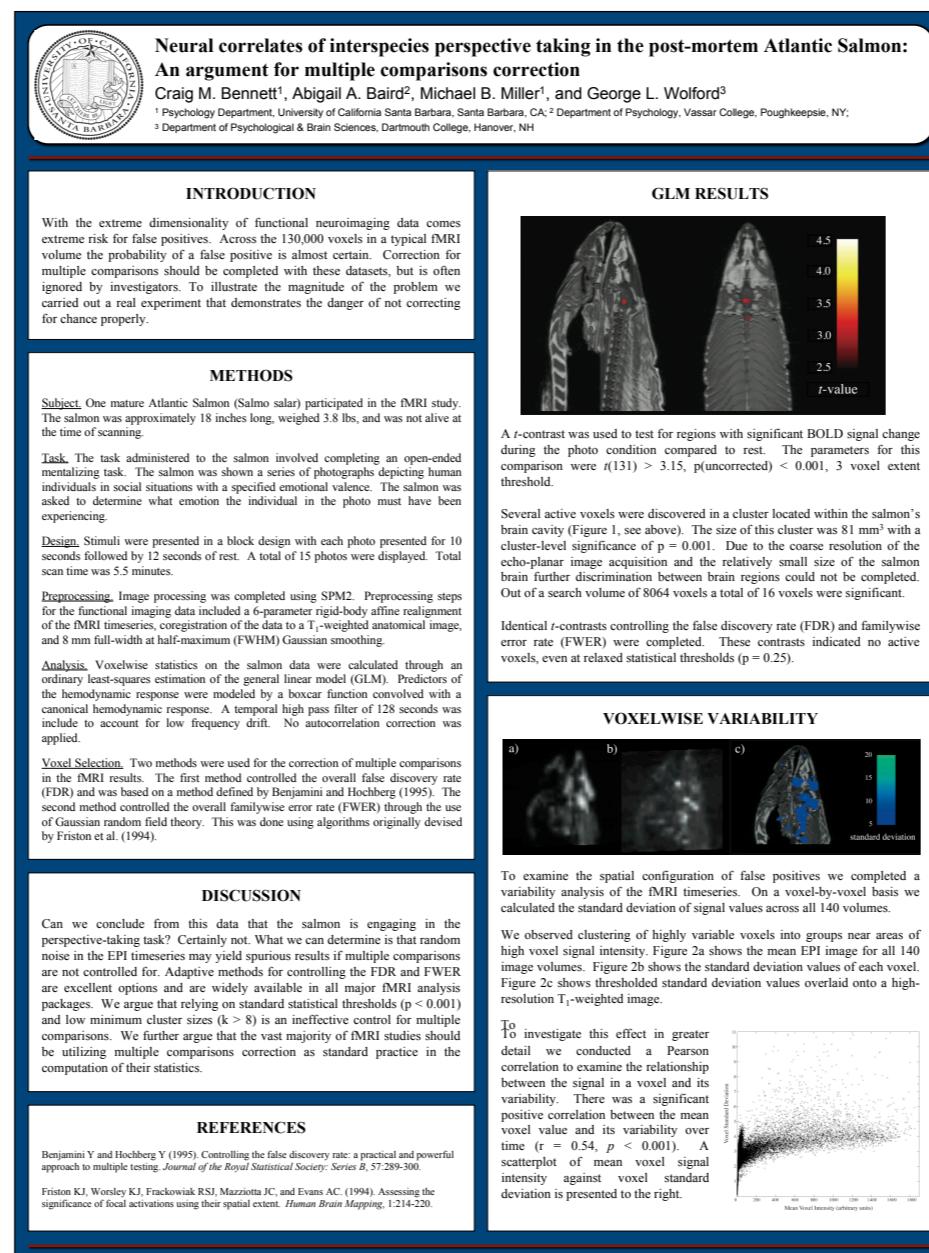
Data “Science”

- To be fair...
 - Intuition plays a huge role in the scientific method (“make observations” is Step 1).
 - Exploratory analysis is necessary, it's okay to not be all rigor all the time
- But!
 - Exploratory analysis (even when it involves the biggest of data) is meant to **form** a hypothesis, not test one
 - Good experimental design and rigorous statistics are essential if we want to make claims about how the world works

Data “Science”

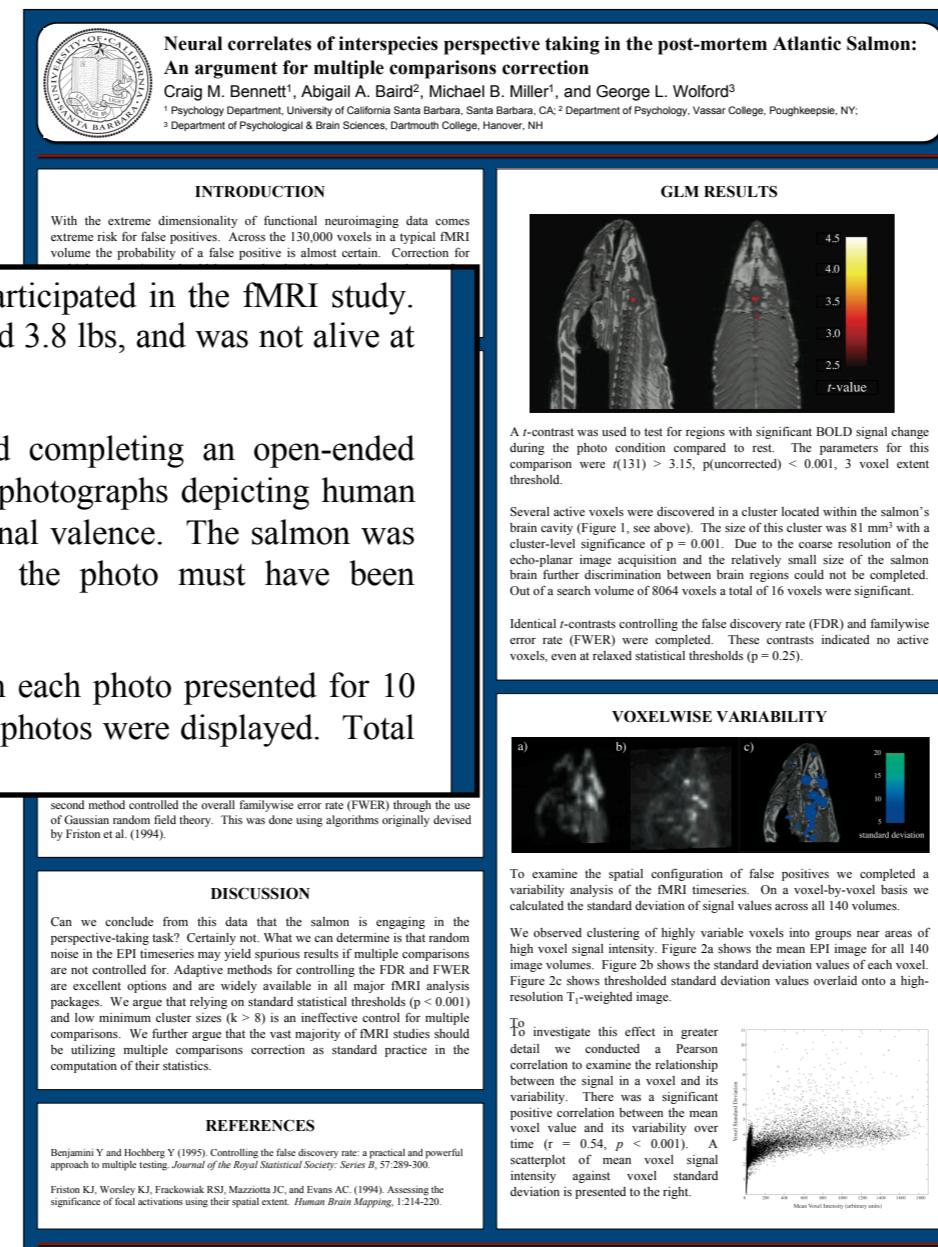


Data “Science”



Neural correlates of interspecies perspective taking in the post-mortem Atlantic Salmon

Data “Science”



Neural correlates of interspecies perspective taking in the post-mortem Atlantic Salmon

Data “Science”

**Neural correlates of interspecies perspective taking in the post-mortem Atlantic Salmon:
An argument for multiple comparisons correction**

Craig M. Bennett¹, Abigail A. Baird², Michael B. Miller¹, and George L. Wolford³

¹ Psychology Department, University of California Santa Barbara, Santa Barbara, CA; ² Department of Psychology, Vassar College, Poughkeepsie, NY;
³ Department of Psychological & Brain Sciences, Dartmouth College, Hanover, NH

INTRODUCTION

With the extreme dimensionality of functional neuroimaging data comes extreme risk for false positives. Across the 130,000 voxels in a typical fMRI volume the probability of a false positive is almost certain. Correction for multiple comparisons should be completed with these datasets, but is often ignored by investigators. To illustrate the magnitude of the problem we carried out a real experiment that demonstrates the danger of not correcting for chance properly.

METHODS

Subject. One mature Atlantic Salmon (*Salmo salar*) participated in the fMRI study. The salmon was approximately 18 inches long, weighed 3.8 lbs, and was not alive at the time of scanning.

Task. The task administered to the salmon involved completing an open-ended mentalizing task. The salmon was shown a series of photographs depicting human individuals in social situations with a specified emotional valence. The salmon was asked to determine what emotion the individual in the photo must have been experiencing.

Design. Stimuli were presented in a block design with each photo presented for 10 seconds followed by 12 seconds of rest. A total of 15 photos were displayed. Total scan time was 5.5 minutes.

Preprocessing. Image processing was completed using SPM2. Preprocessing steps for the functional imaging data included a 6-parameter rigid-body affine realignment of the fMRI timeseries, coregistration of the data to a T_1 -weighted anatomical image, and 8 mm full-width at half-maximum (FWHM) Gaussian smoothing.

Analysis. Voxelwise statistics on the salmon data were calculated through an ordinary least-squares estimation of the general linear model (GLM). Predictors of the hemodynamic response were modeled by a boxcar function convolved with a canonical hemodynamic response. A temporal high pass filter of 128 seconds was applied to account for low frequency drift. No autocorrelation correction was applied.

Voxel Selection. Two methods were used for the correction of multiple comparisons in the fMRI results. The first method controlled the overall false discovery rate (FDR) and was based on a method defined by Benjamini and Hochberg (1995). The second method controlled the overall familywise error rate (FWER) through the use of Gaussian random field theory. This was done using algorithms originally devised by Friston et al. (1994).

DISCUSSION

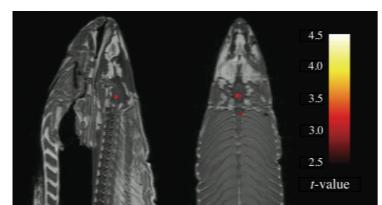
Can we conclude from this data that the salmon is engaging in the perspective-taking task? Certainly not. What we can determine is that random noise in the EPI timeseries may yield spurious results if multiple comparisons are not controlled for. Adaptive methods for controlling the FDR and FWER are excellent options and are widely available in all major fMRI analysis packages. We argue that relying on standard statistical thresholds ($p < 0.001$) and low minimum cluster sizes ($k > 8$) is an ineffective control for multiple comparisons. We further argue that the vast majority of fMRI studies should be utilizing multiple comparisons correction as standard practice in the computation of their statistics.

REFERENCES

Benjamini Y and Hochberg Y (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society: Series B*, 57, 289-300.

Friston KJ, Worsley KJ, Frackowiak RSJ, Mazziotta JC, and Evans AC. (1994). Assessing the significance of focal activations using their spatial extent. *Human Brain Mapping*, 1:214-220.

GLM RESULTS



A t -contrast was used to test for regions with significant BOLD signal change during the photo condition compared to rest. The parameters for this comparison were $t(131) > 3.15$, $p(\text{uncorrected}) < 0.001$, 3 voxel extent threshold.

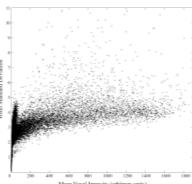
Several active voxels were discovered in a cluster located within the salmon's brain cavity (Figure 1, see above). The size of this cluster was 81 mm^3 with a cluster-level significance of $p = 0.001$. Due to the coarse resolution of the echo-planar image acquisition and the relatively small size of the salmon brain further discrimination between brain regions could not be completed. Out of a search volume of 8064 voxels a total of 16 voxels were significant.

Identical t -contrasts controlling the false discovery rate (FDR) and familywise error rate (FWER) were completed. These contrasts indicated no active voxels, even at relaxed statistical thresholds ($p = 0.25$).

Can we conclude from this data that the salmon is engaging in the perspective-taking task? Certainly not. What we can determine is that random noise in the EPI timeseries may yield spurious results if multiple comparisons are not controlled for. Adaptive methods for controlling the FDR and FWER

We observed clustering of highly variable voxels into groups near areas of high voxel signal intensity. Figure 2a shows the mean EPI image for all 140 image volumes. Figure 2b shows the standard deviation values of each voxel. Figure 2c shows thresholded standard deviation values overlaid onto a high-resolution T_1 -weighted image.

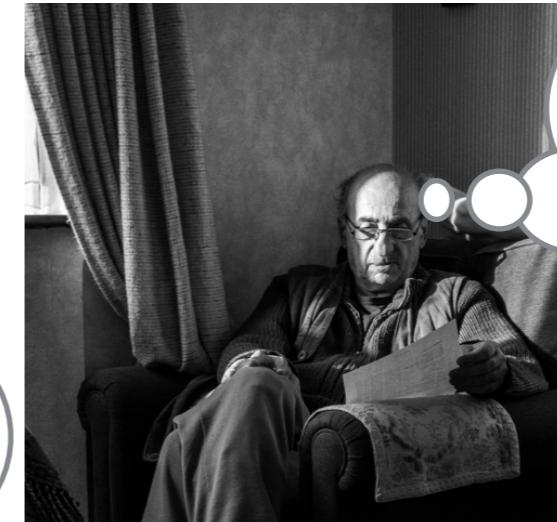
To investigate this effect in greater detail we conducted a Pearson correlation to examine the relationship between the signal in a voxel and its variability. There was a significant positive correlation between the mean voxel value and its variability over time ($r = 0.54$, $p < 0.001$). A scatterplot of mean voxel signal intensity against voxel standard deviation is presented to the right.



Neural correlates of interspecies perspective taking in the post-mortem Atlantic Salmon

“Data” Science

“Data” Science





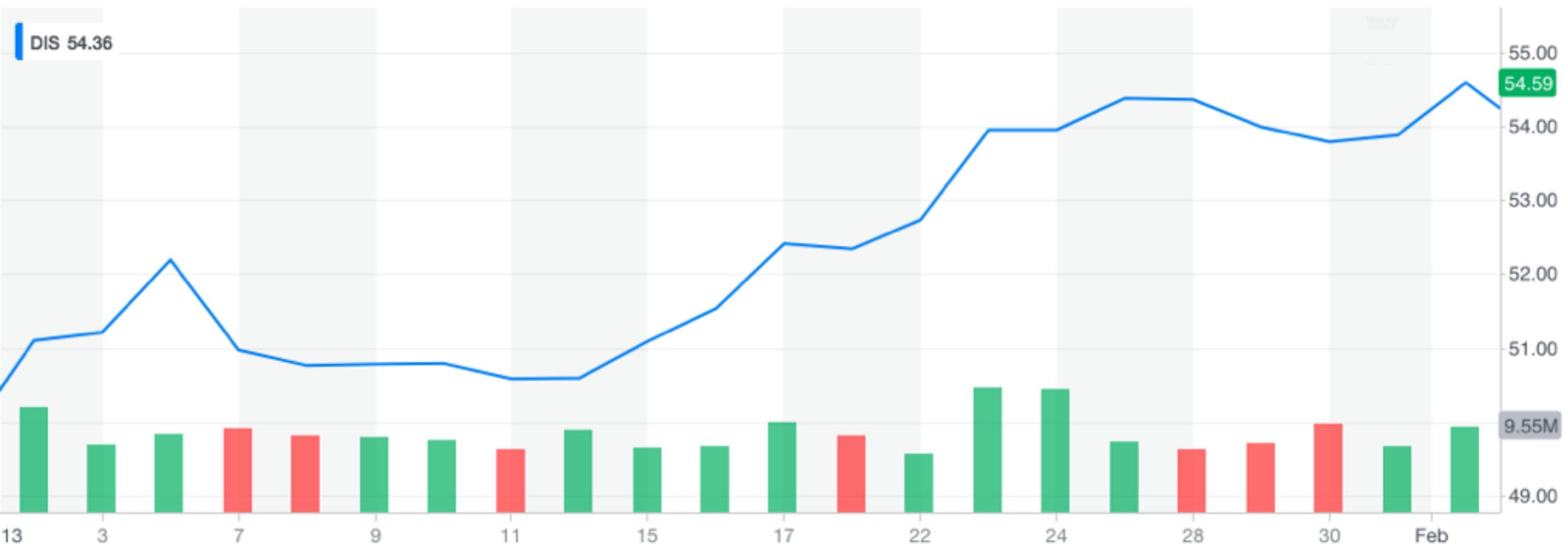
Roses are red.
Violets are blue.





Roses are red
Violets are blue

“Data” Science



“Data” Science



“Data” Science



“Data” Science



“Data” Science

- To be fair...

“Data” Science

- To be fair...
 - Not all science is empirical—its possible to gain insight and make progress via introspection

“Data” Science

- To be fair...
 - Not all science is empirical—its possible to gain insight and make progress via introspection
 - E.g. simulations, case studies, motivating/illustrative examples, worst-case vs. average case runtime

“Data” Science

- To be fair...
 - Not all science is empirical—its possible to gain insight and make progress via introspection
 - E.g. simulations, case studies, motivating/illustrative examples, worst-case vs. average case runtime
- But!

“Data” Science

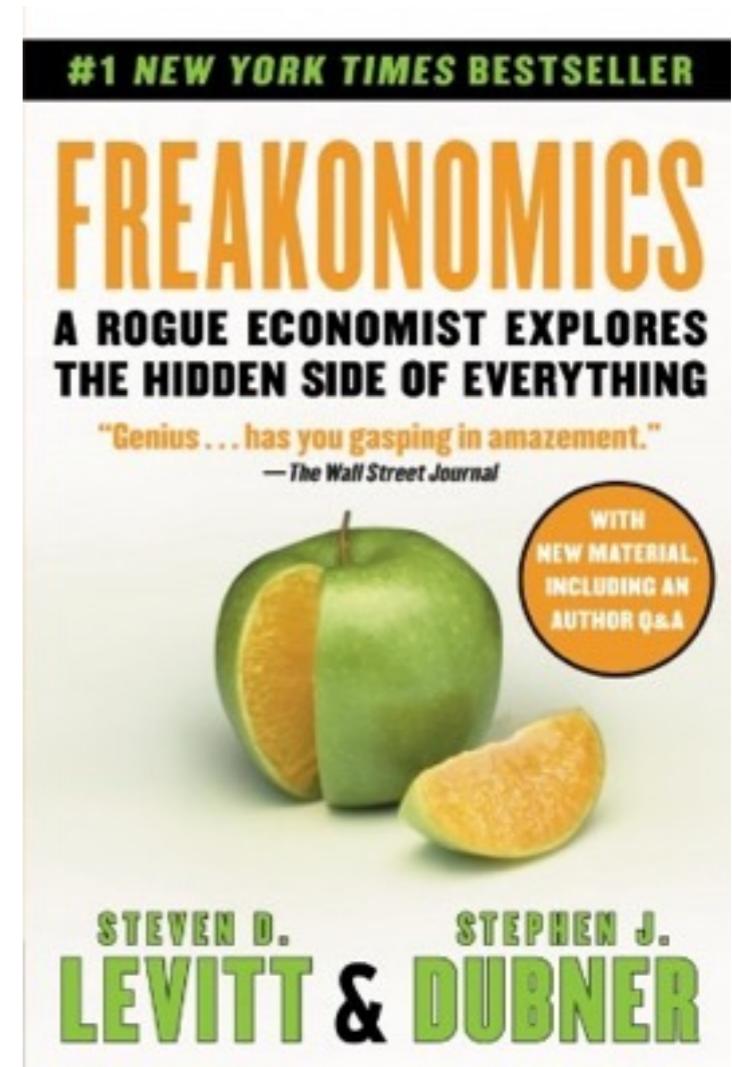
- To be fair...
 - Not all science is empirical—its possible to gain insight and make progress via introspection
 - E.g. simulations, case studies, motivating/illustrative examples, worst-case vs. average case runtime
- But!
 - Theory is only helpful if it mirrors practice.

“Data” Science

- To be fair...
 - Not all science is empirical—its possible to gain insight and make progress via introspection
 - E.g. simulations, case studies, motivating/illustrative examples, worst-case vs. average case runtime
- But!
 - Theory is only helpful if it mirrors practice.
 - “All models are wrong, but some are useful.”

“Data” Science

- Problem: Parents run late when picking kids up from day care
- Sensible Solution: Impose a late fee



<https://www.nytimes.com/2005/05/15/books/chapters/freakonomics.html>

<https://rady.ucsd.edu/faculty/directory/gneezy/pub/docs/fine.pdf>

“Data” Science

- Problem: Parents run late when picking kids up from day care
- Sensible Solution: Impose a late fee

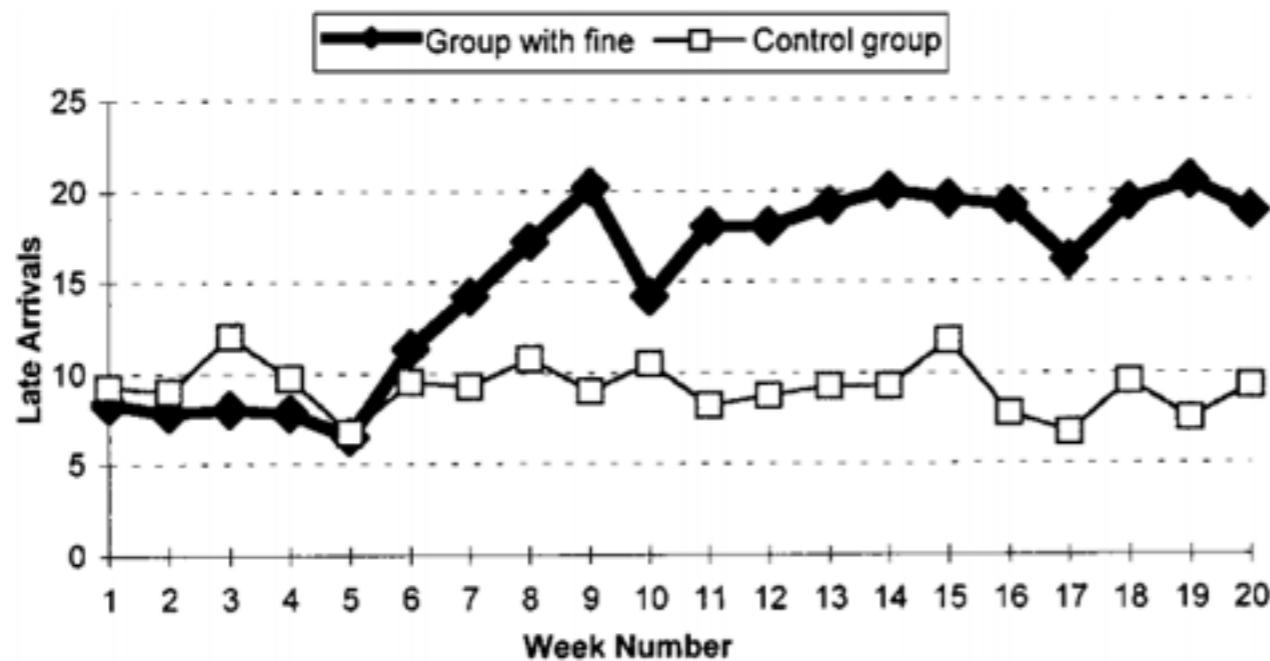
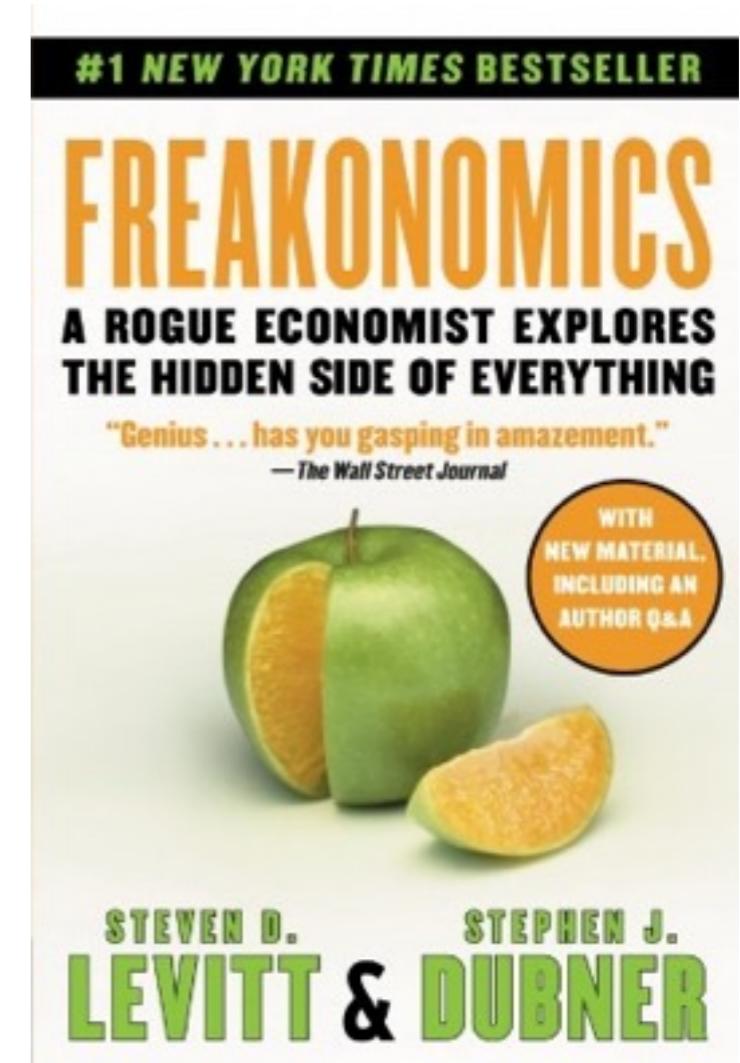


FIGURE 1.—Average number of late-coming parents, per week



<https://www.nytimes.com/2005/05/15/books/chapters/freakonomics.html>

<https://rady.ucsd.edu/faculty/directory/gneezy/pub/docs/fine.pdf>

“Data” Science

- He is not lucky to have to pay for the property.



“Data” Science

- He is not lucky to have to pay for the property.
 - Did he pay for the property?

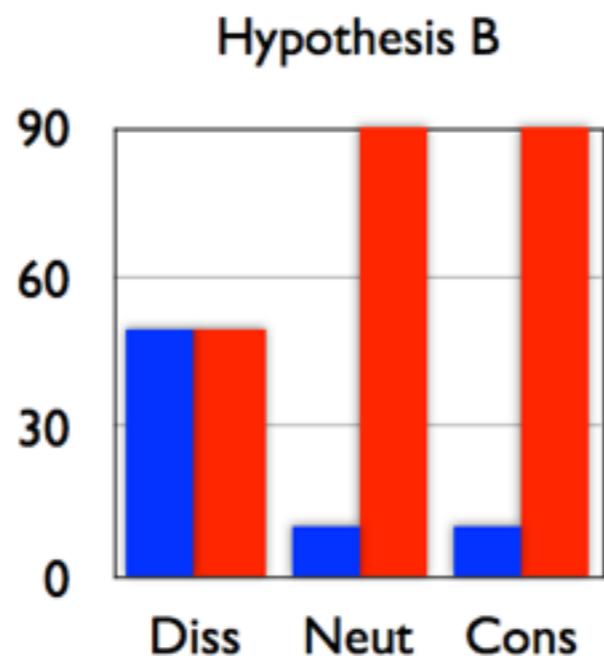
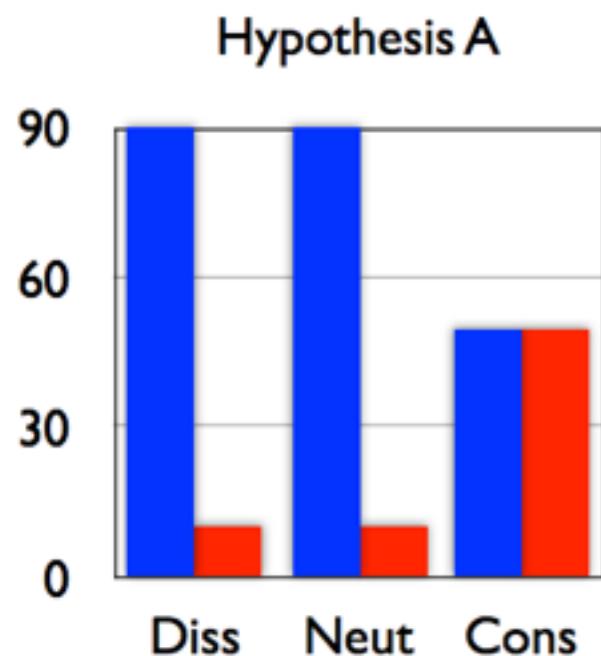


“Data” Science

- He is not lucky to have to pay for the property.
 - Did he pay for the property?
- The girl was not lucky to get away alive.
 - Did she get away alive?



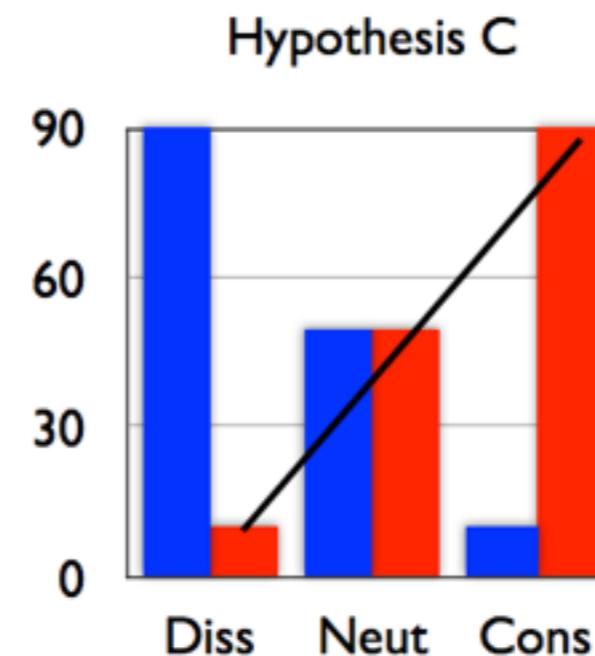
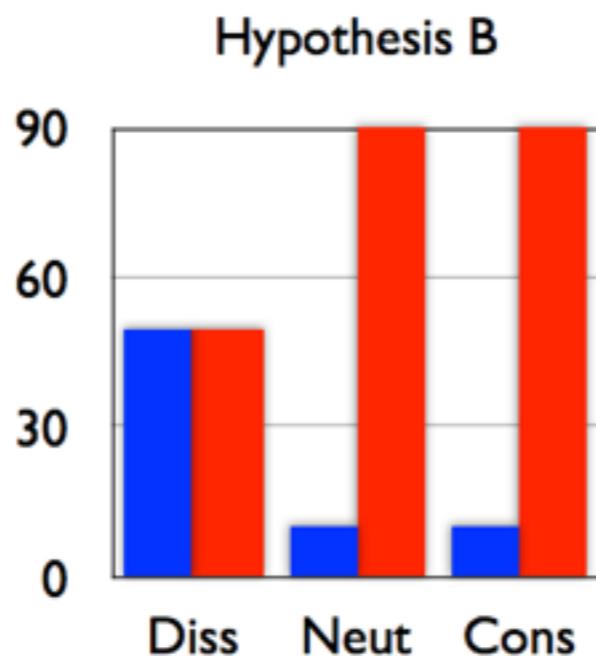
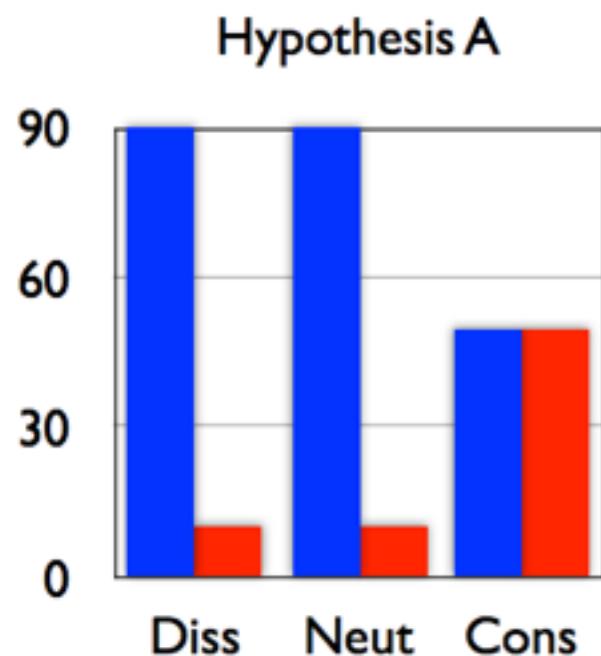
“Data” Science



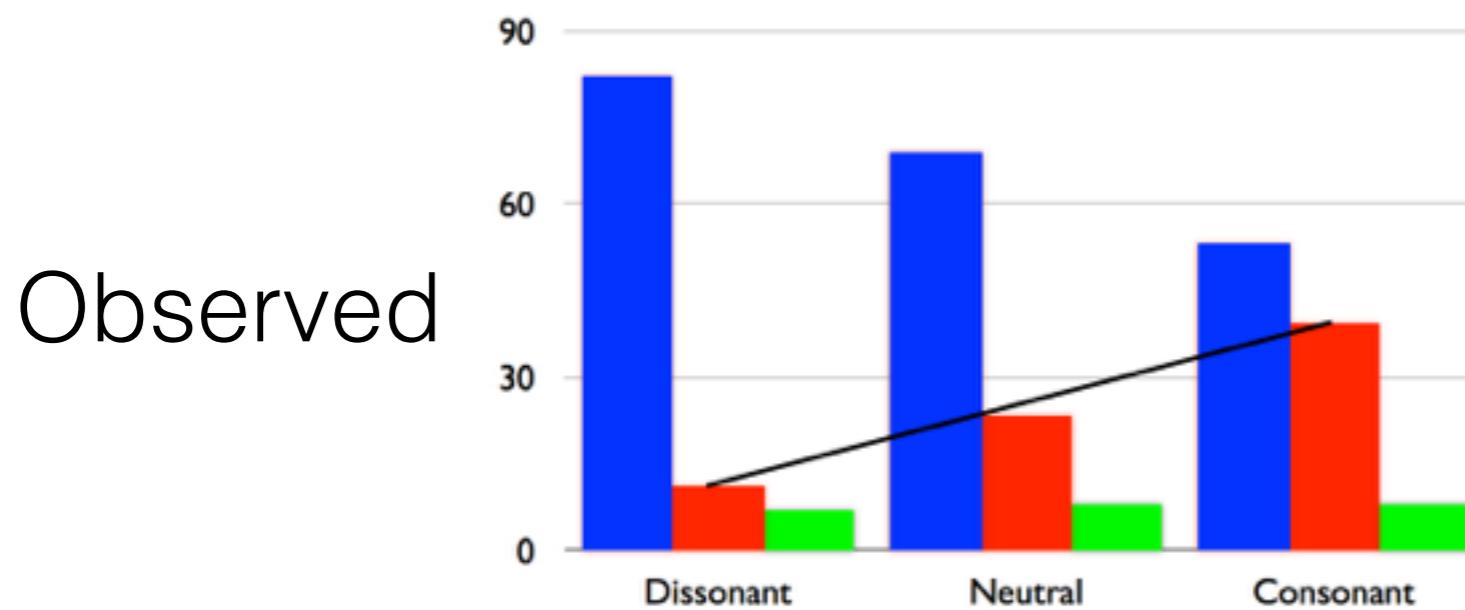
YES NO



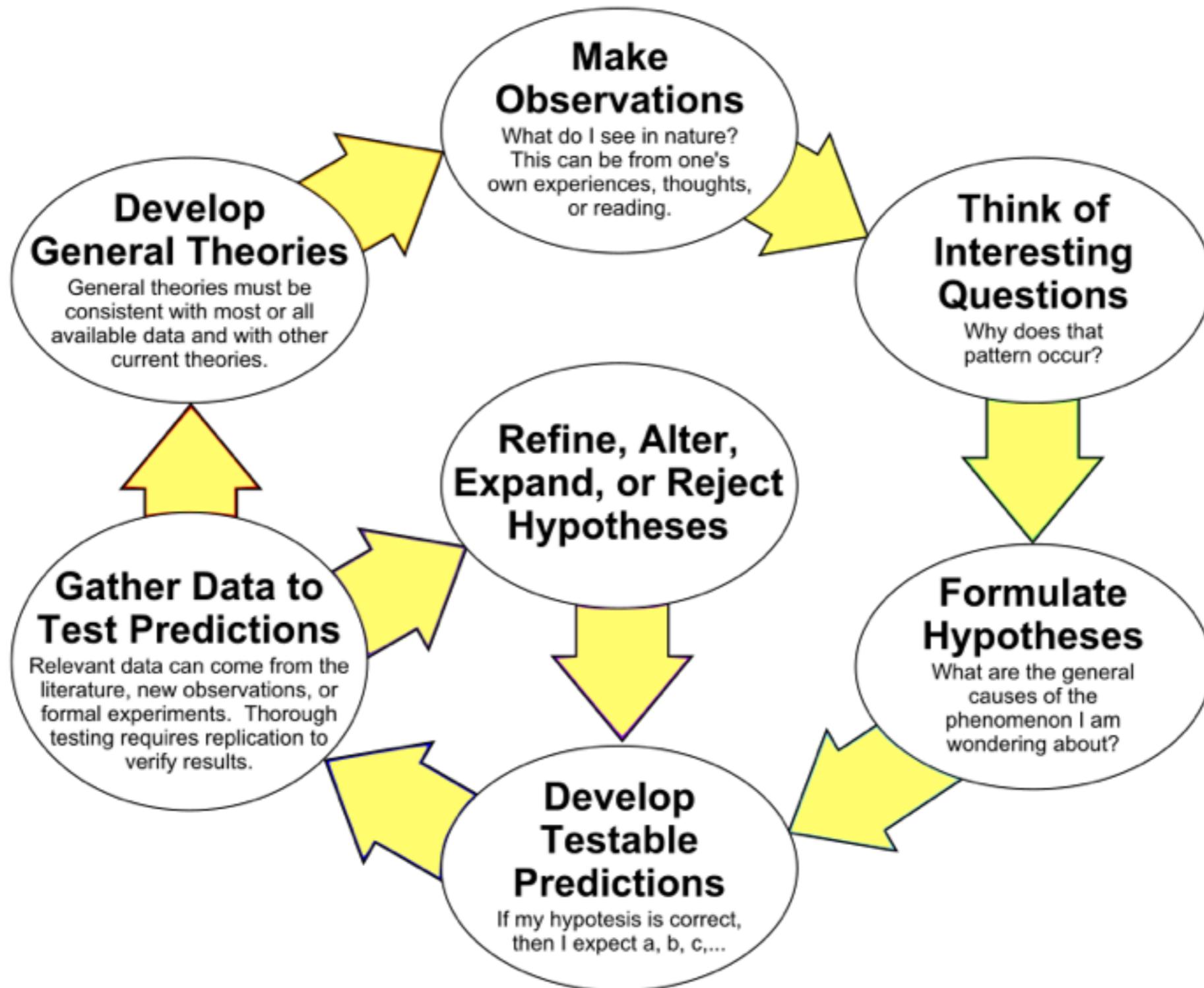
“Data” Science



YES NO



Data! Science!



CSCI 19S1A

What is ~~Data Science?~~

January

S	M	T	W	T	F	S
			1	2	3	4
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31		

February

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28		

March

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30
31						

April

S	M	T	W	T	F	S
		2	3	4	5	6
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30				

- Data Collection/Cleaning

- Probability and Statistics

- Machine Learning

- Advanced Topics/
Applications

- Other Topics

January

S	M	T	W	T	F	S
			1	2	3	4
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31		

February

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28		

March

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30
31						

April

S	M	T	W	T	F	S
			1	2	3	4
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30				

Right Here, Right Now.

January

S	M	T	W	T	F	S
			1	2	3	4
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31		

February

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28		

March

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30
31						

April

S	M	T	W	T	F	S
			1	2	3	4
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30				

Databases for Data Scientists:

- Entity-Relationship (ER) Diagrams
- Relational Algebra
- SQL
- [Briefly] Optimization
- [Briefly] NoSQL

January

S	M	T	W	T	F	S
			1	2	3	4
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31		

Collecting and Cleaning Data:

February

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28		

March

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30
31						

April

S	M	T	W	T	F	S
			1	2	3	4
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30				

- Crowdsourcing

- Cleaning, Normalization,
Regular Expressions

- Web Crawling

- APIs

January

S	M	T	W	T	F	S
			1	2	3	4
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31		

February

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28		

March

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30
31						

April

S	M	T	W	T	F	S
			1	2	3	4
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30				

Big Data and Working at Scale

- Massively Parallel Processing (MapReduce, Storm)
- [Briefly] Randomized Data Structures

January

S	M	T	W	T	F	S
			1	2	3	4
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31		

February

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28		

March

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30
31						

April

S	M	T	W	T	F	S
			1	2	3	4
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30				

Intro to Probability

- Random Variables
- Sample Spaces
- Distributions
- Notation

January

S	M	T	W	T	F	S
			1	2	3	4
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31		

February

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28		

March

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30
31						

April

S	M	T	W	T	F	S
			1	2	3	4
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30				

Hypothesis Testing

- Central Limit Theorem
- P-Values
- T-Tests, Chi-Squared Tests
- Regression
- Fixed and Random Effects

January

S	M	T	W	T	F	S
			1	2	3	4
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31		

Intro to ML

February

S	M	T	W	T	F	S
				1	2	
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28		

March

S	M	T	W	T	F	S
				1	2	
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30
31						

April

S	M	T	W	T	F	S
			1	2	3	4
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30				

- Feature Representations
- Loss Functions
- Supervised vs. Unsupervised Learning
- Overview of Categorizations of Models

January

S	M	T	W	T	F	S
			1	2	3	4
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31		

February

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28		

March

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30
31						

April

S	M	T	W	T	F	S
			1	2	3	4
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30				

ML for Data Scientists

- Clustering and Nearest Neighbors
- Linear Regression, Logistic Regression, and SVMs
- Estimating Parameters with Gradient Descent
- Using SciKit Learn

January

S	M	T	W	T	F	S
			1	2	3	4
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31		

February

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28		

March

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30
31						

April

S	M	T	W	T	F	S
			1	2	3	4
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30				

Trouble Shooting ML

- Overfitting and Generalizability
- Regularization
- Feature Selection

January

S	M	T	W	T	F	S
			1	2	3	4
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31		

February

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28		

March

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30
31						

April

S	M	T	W	T	F	S
			1	2	3	4
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30				

Visualization

- Best Practices
- Using D3 and matplotlib

January

S	M	T	W	T	F	S
			1	2	3	4
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31		

February

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28		

March

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30
31						

April

S	M	T	W	T	F	S
		1	2	3	4	5
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30				

Trouble Shooting ML

- Sampling
- Evaluation Metrics

January

S	M	T	W	T	F	S
			1	2	3	4
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31		

February

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28		

March

S	M	T	W	T	F	S
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30
31						

April

S	M	T	W	T	F	S
			1	2	3	4
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30				

How to Lie with Statistics

- p-hacking
- Researcher Degrees of Freedom
- Issues with Reproducibility

January

S	M	T	W	T	F	S
			1	2	3	4
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31		

Special Topics/Applications

February

S	M	T	W	T	F	S
				1	2	
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28		

March

S	M	T	W	T	F	S
				1	2	
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30
31						

April

S	M	T	W	T	F	S
		1	2	3	4	5
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30				

- NLP, Topic Modeling
- Algorithmic Bias, Ethics
- Recommendation Systems
- Deep Learning
- Causal Inference

Assignments

- Feb 9: SQL and Query Optimization
- Feb 14: Web Crawling and Data Cleaning
- Feb 28: Map Reduce
- Mar 7: Linear Regression
- Mar 21: K Means
- Apr 8: Visualization
- Apr 18: Topic Modeling
- Apr 25: Deep Learning

Project

- Feb 18: Pre-Proposal Due (10%)
- Mar 8: Check-in 1
- Mar 15: Blog Post 1 (10%)
- Apr 5: Midterm Report (30%)
- Apr 16: Check-in 2
- Apr 19: Blog Post 2 (10%)
- May 3: Posters Due
- May 6-7: Poster Presentation (20%)
- May 10: Blog Post 3 (Final Writeup) (20%)

Grading

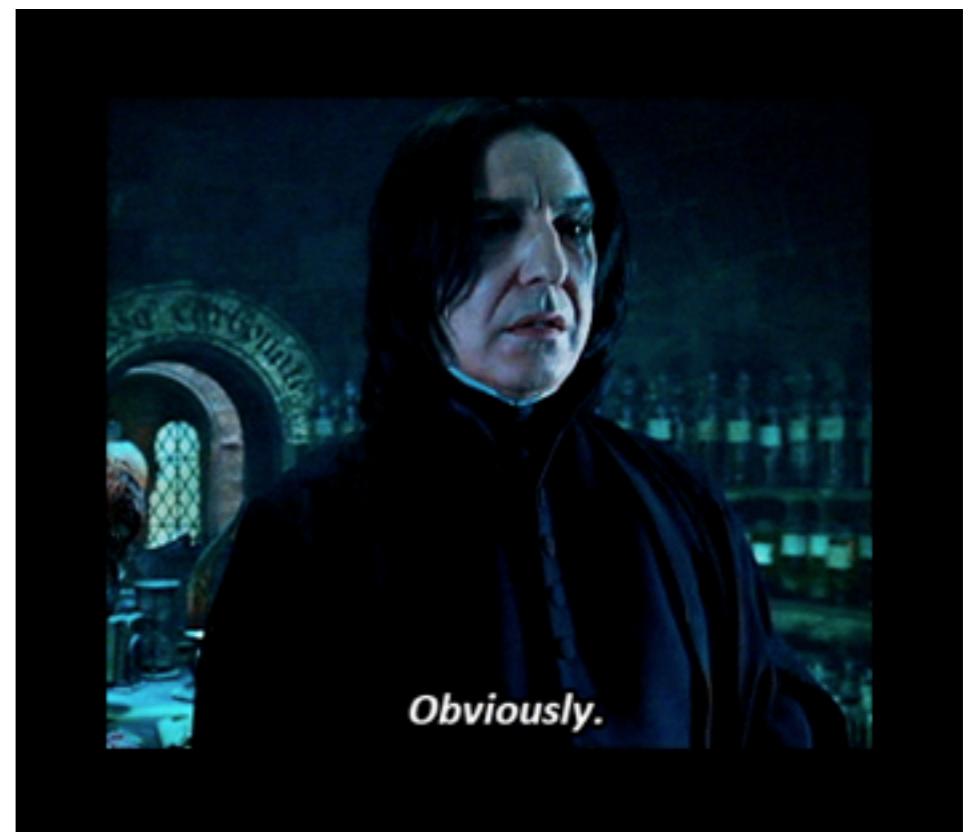
- 60% Assignments (7.5% each)
- 30% Final Project
- 10% Attendance/Clickers (must attend 2/3 of classes)

Late Days

- Assignments are due at 11:59 pm on the listed due date
- 5 late days total; maximum of 2 on any single assignment
- 20% penalty for each additional day late
- No late days for Final Project deliverables (incl. intermediate deliverables)

Collaboration

- Talking to each other is good. Cheating is bad.
- Sign the form so I know you know.

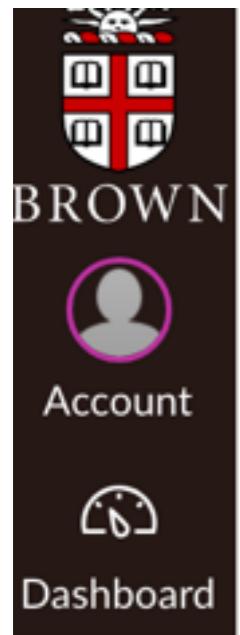


To Do Now

To Do Now

- Get on the waitlist—make your case there. (Please don't send emails to me directly.)

To Do Now



≡ Spring 2019 CSCI

2019 Spring

Home

Discussions

Grades

People

Syllabus

Media Library

Collaborations

Chat

iClicker Sync

- Join iClicker: <https://ithelp.brown.edu/kb/articles/iclicker-cloud-reef-instructions-for-students>
- Make sure you register via canvas so that grades get synced

To Do Now

- Join the course on Piazza
- Piazza is now opt-out (as opposed to opt-in) for data sharing.
- Decide how you feel about this. Instructions for opt-out are on Canvas.

2019 Spring

PAGE TITLE ▾

[Home](#)

[Piazza and Student Privacy - 2019](#)

[Discussions](#)

[Grades](#)

[People](#)

[Pages](#)

[Files](#)

To Do Now

- Hours are starting this week! Go say hi to your staff...



Gurnaaaz



Maulik



Wennie



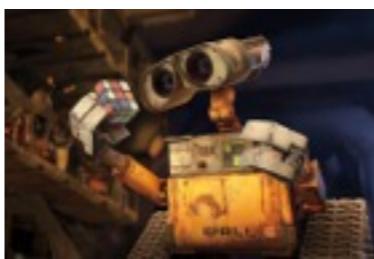
Alex



Ashish



Paarul



Jens



Yiquan



Shivani



Mounika



Fumeng



Pavlo



Tanvir



David



Shre



Zander



Hyunjoon



Erin



Miles



Haomo



Esteban



Iris



Weiqi



Palak

Your Phenomenal Staff!

Thank you!
Questions?