

# What is Data Science?

January 23, 2020

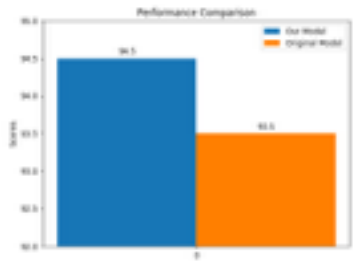
Data Science CSCI 1951A

Brown University

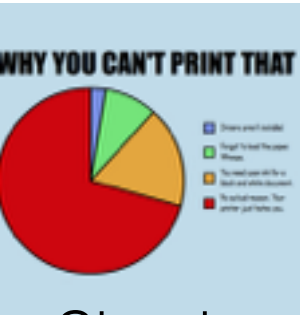
Instructor: Ellie Pavlick

HTAs: Josh Levin, Diane Mutako, Sol Zitter

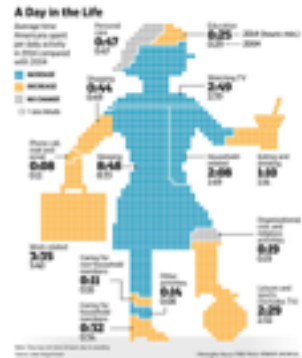
# Your Phenomenal Staff!



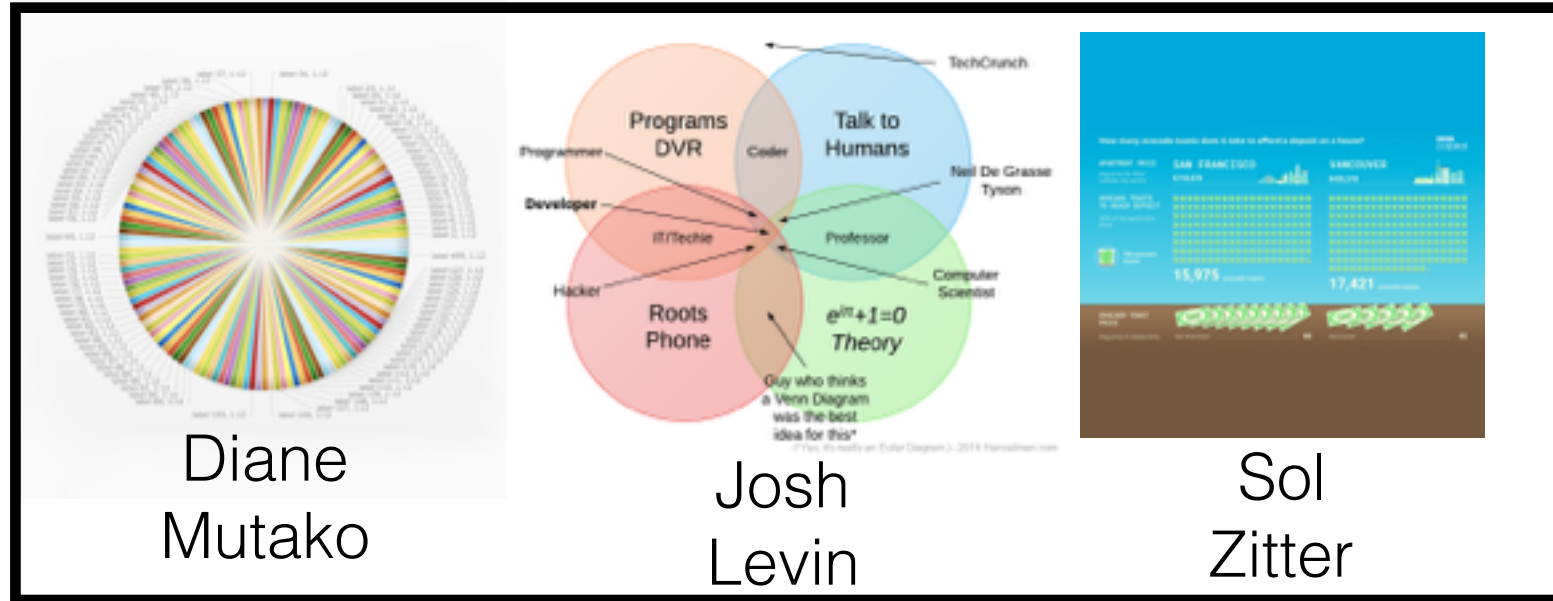
Shunjia Zhu



Shash Sinha



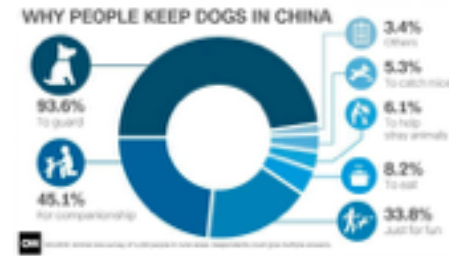
Maggie Wu



Diane Mutako

Josh Levin

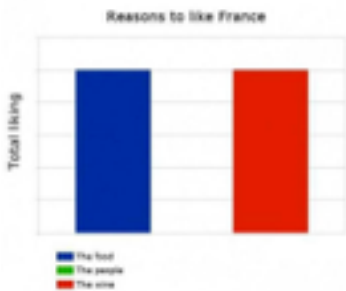
Sol Zitter



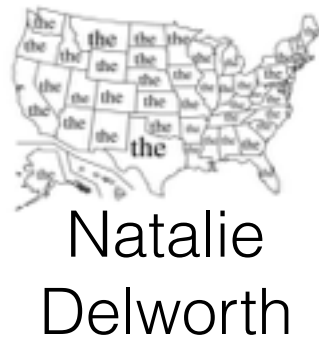
Karlly Feng



Neil Sehgal



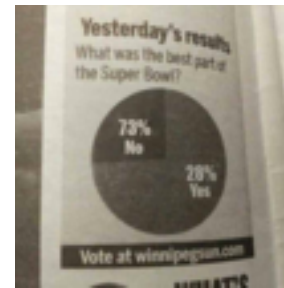
Nazem Aldroubi



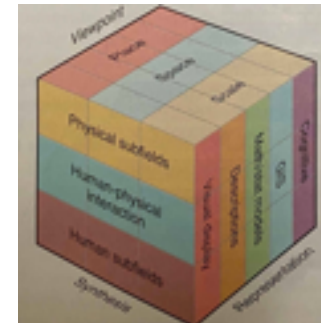
Natalie Delworth



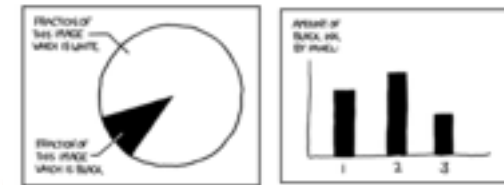
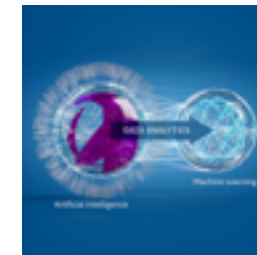
Jonathan Weisskoff



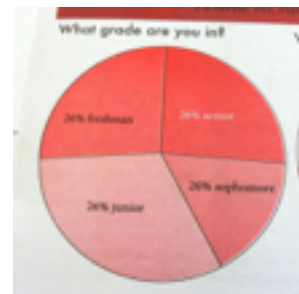
JP Champa



Ben Gershuny



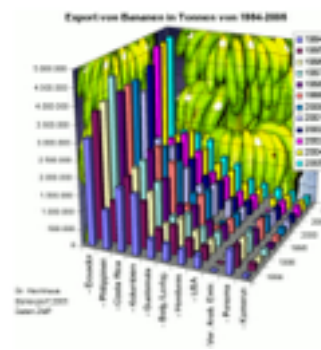
Will Glaser



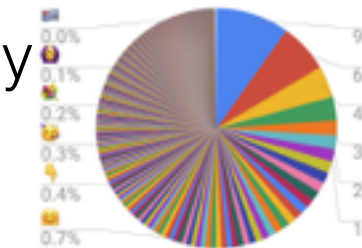
Mounika Dandu



Ben Vu



Marcin Kolaszewski



Sunny Deng



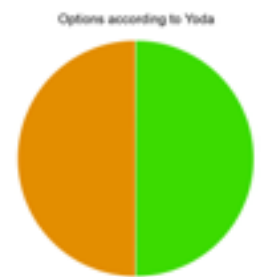
Juho Choi



Arvind Yalavarti



Minna Kimura-



Nam Do

# Waitlist

- If you are not registered, make sure you are on the waitlist (**link is on course webpage**)
- We have a \*little\* wiggle room in the enrollment cap
- We will prioritize fairly (i.e. graduating and need this to graduate > graduating > not graduating...)

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.

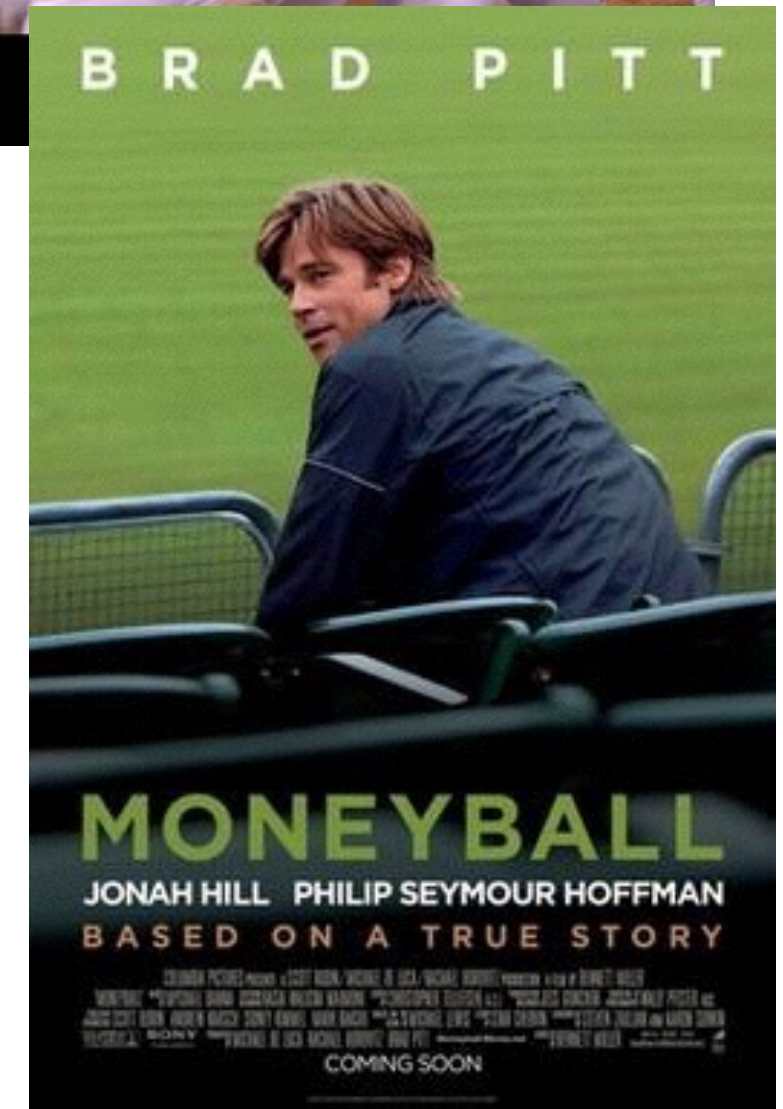
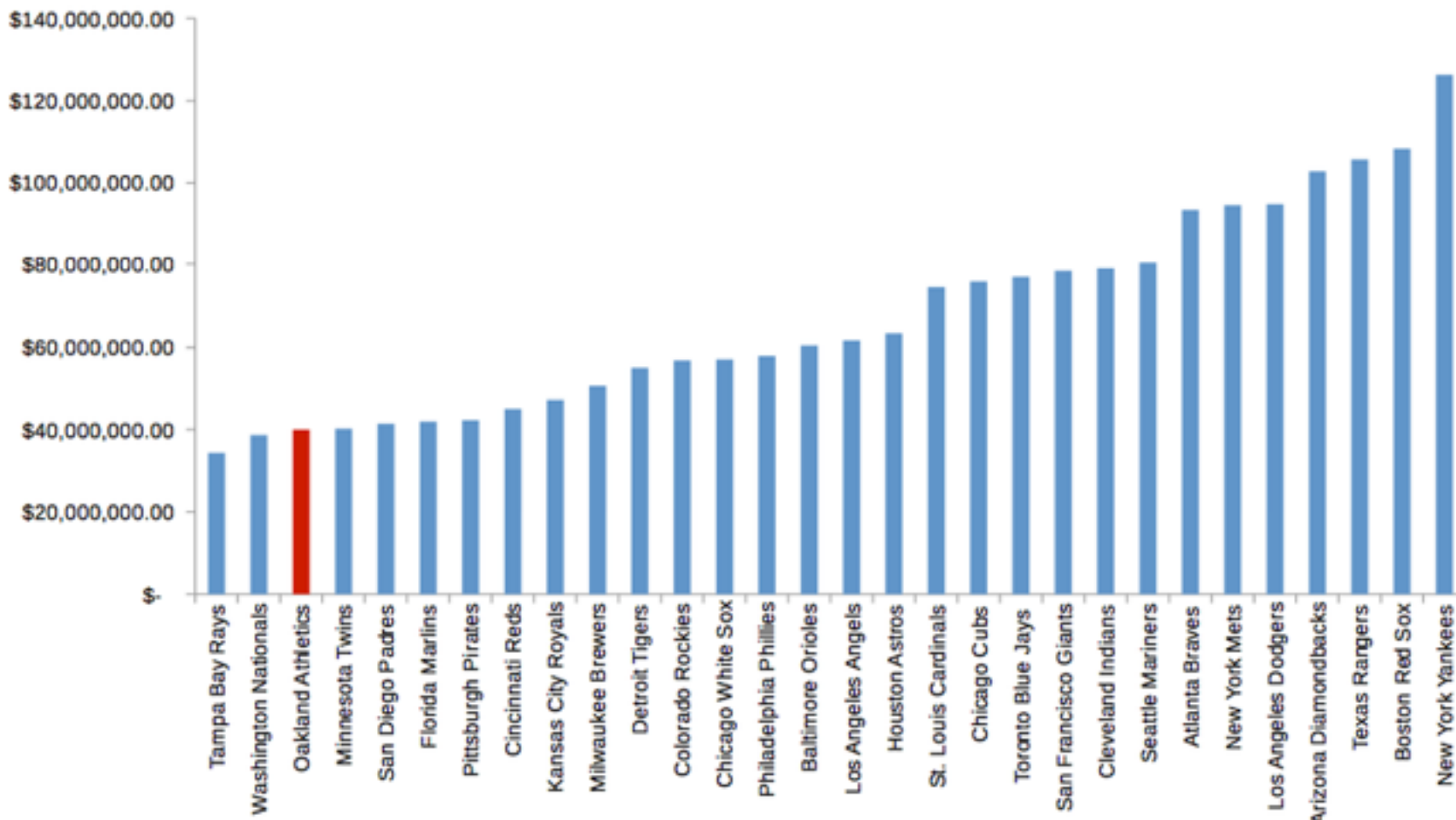
Data scientists realize that they face technical limitations, but they don't allow that to bog down their search for novel solutions. As they make discoveries, they communicate what they've learned and suggest its implications for new business directions. Often they are creative in displaying information 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



↑ +19%

<http://bit.ly/bsack-eg5>



↑ +5%

<http://bit.ly/len>



# Google's “40 Shades of Blue”

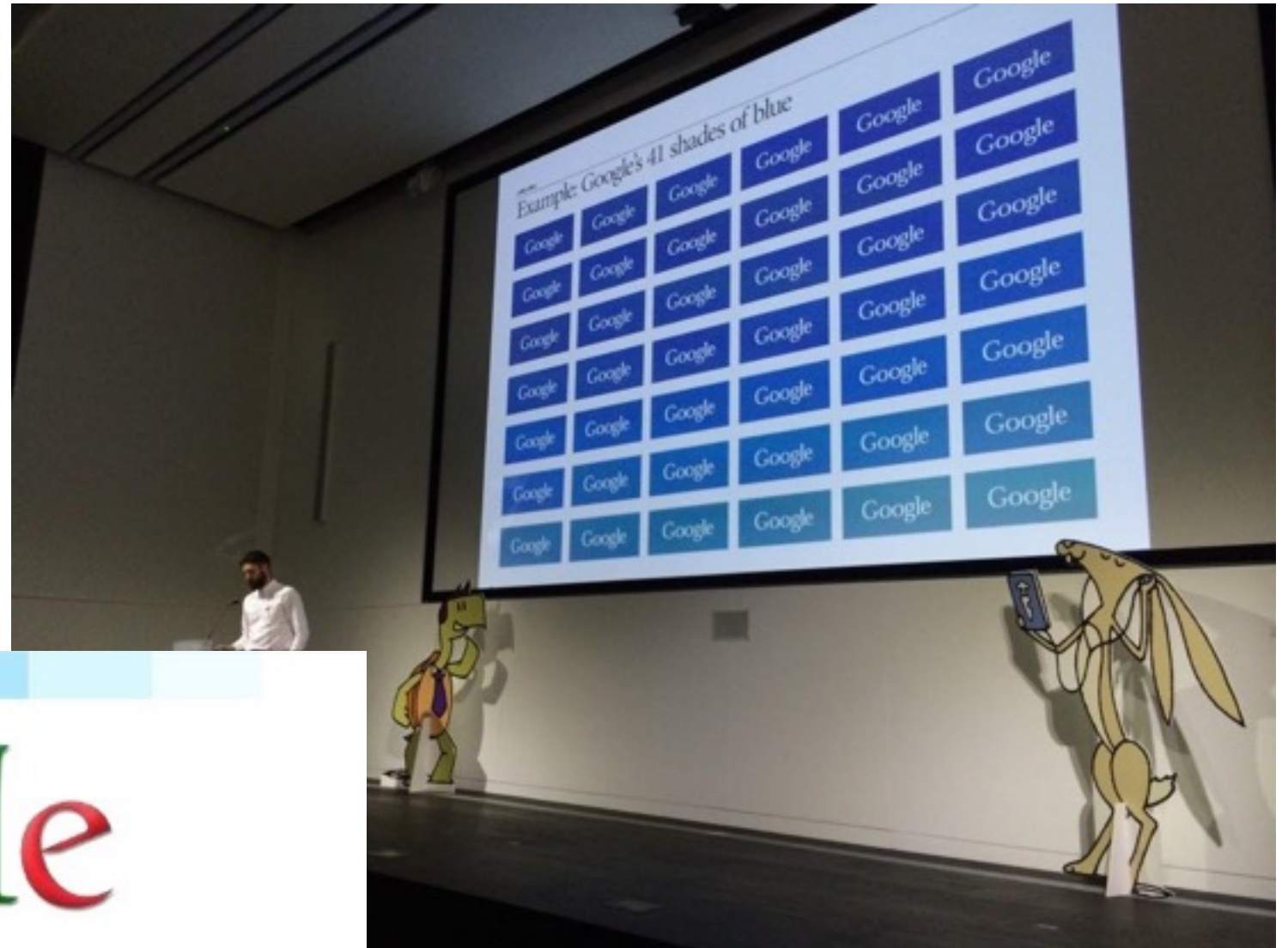


Google

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](https://liveslides.com/download)

**Start the presentation.**



# MACHINE LEARNING

PHOTO/VIDEO  
DATABASE

READING HABITS

CONSUMER  
BEHAVIOR/  
PREFERENCES

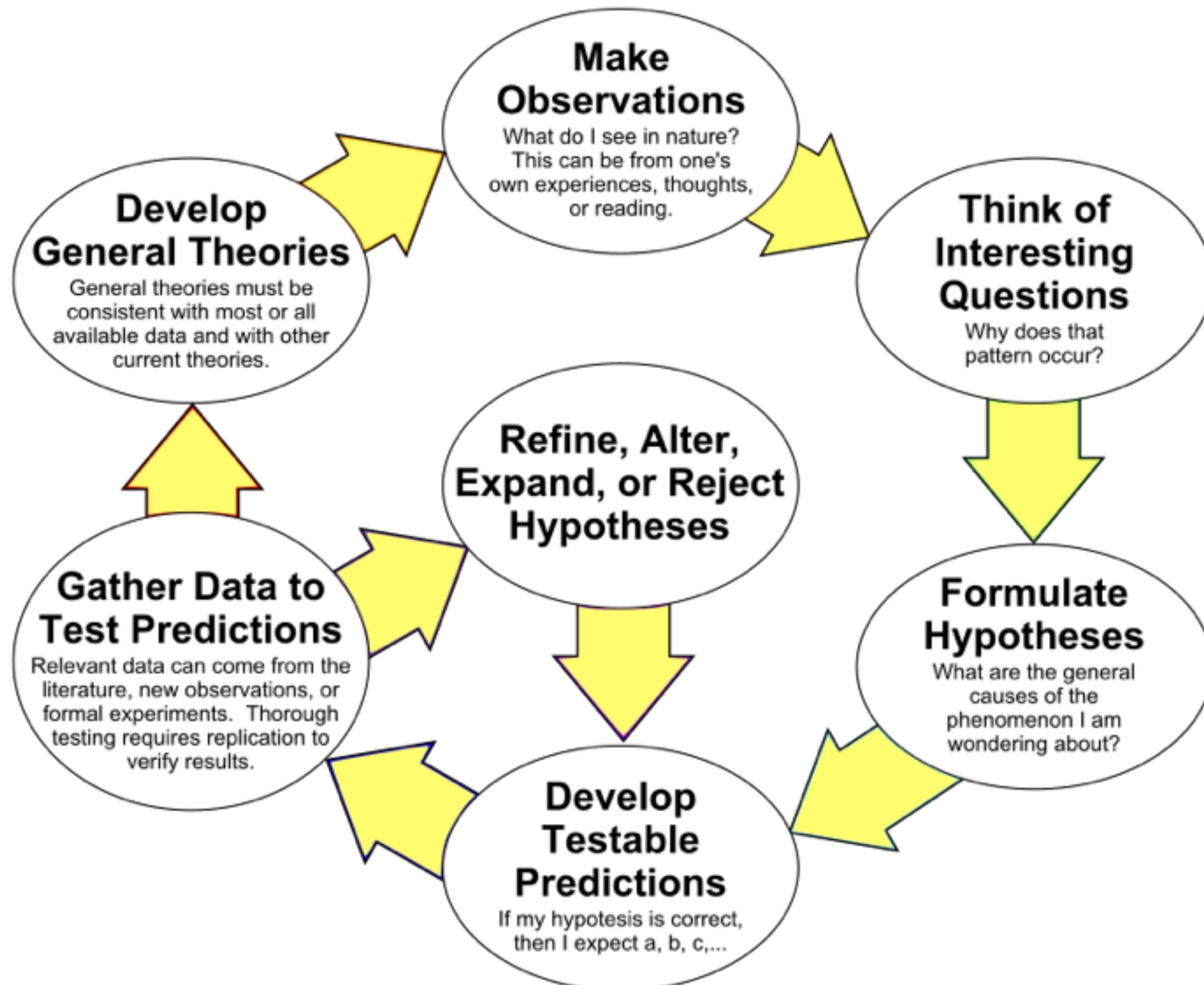


VISUALIZATIONS

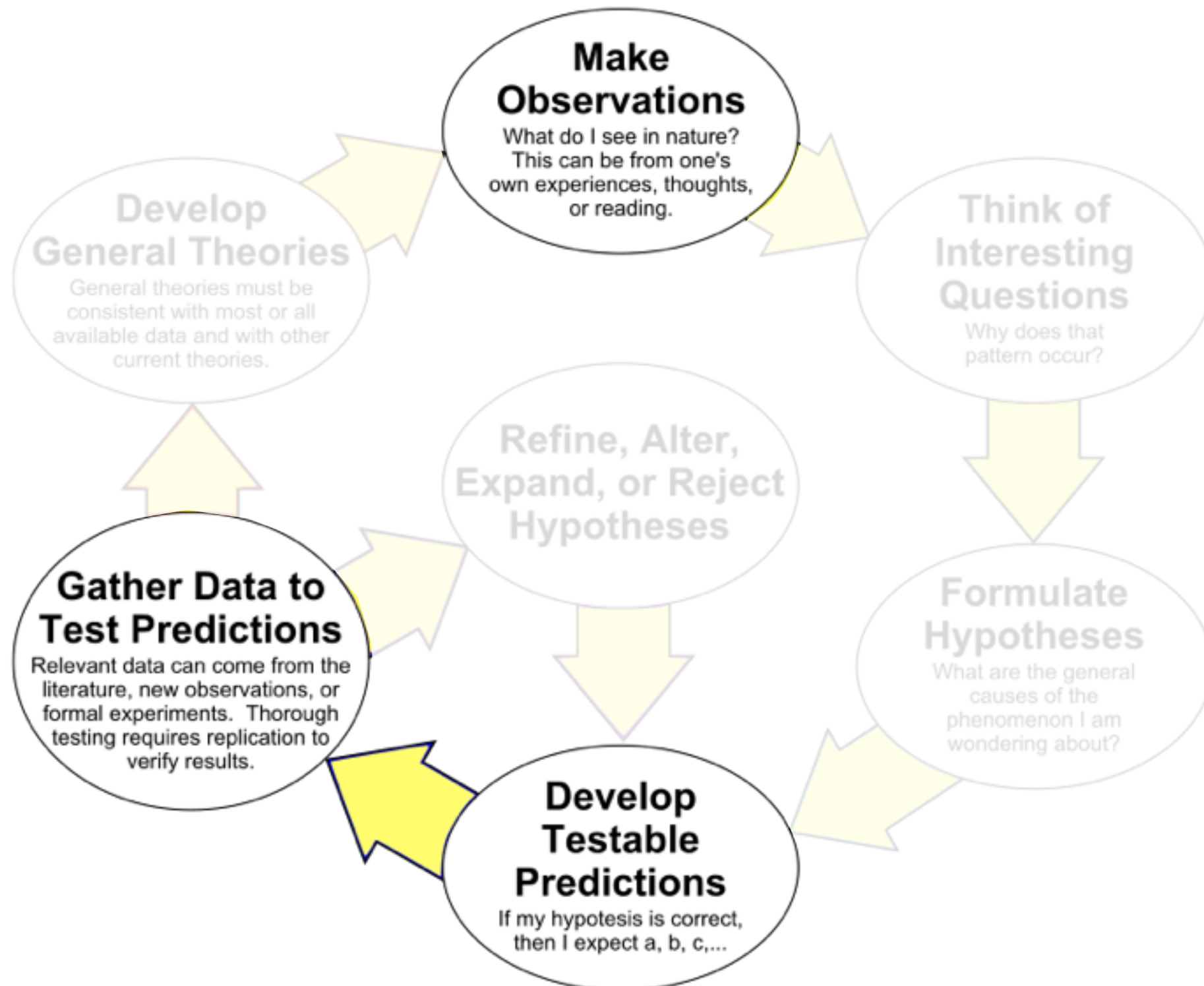
INCREASE CONSUMPTION

HIGH ENGAGEMENT

# The Scientific Method

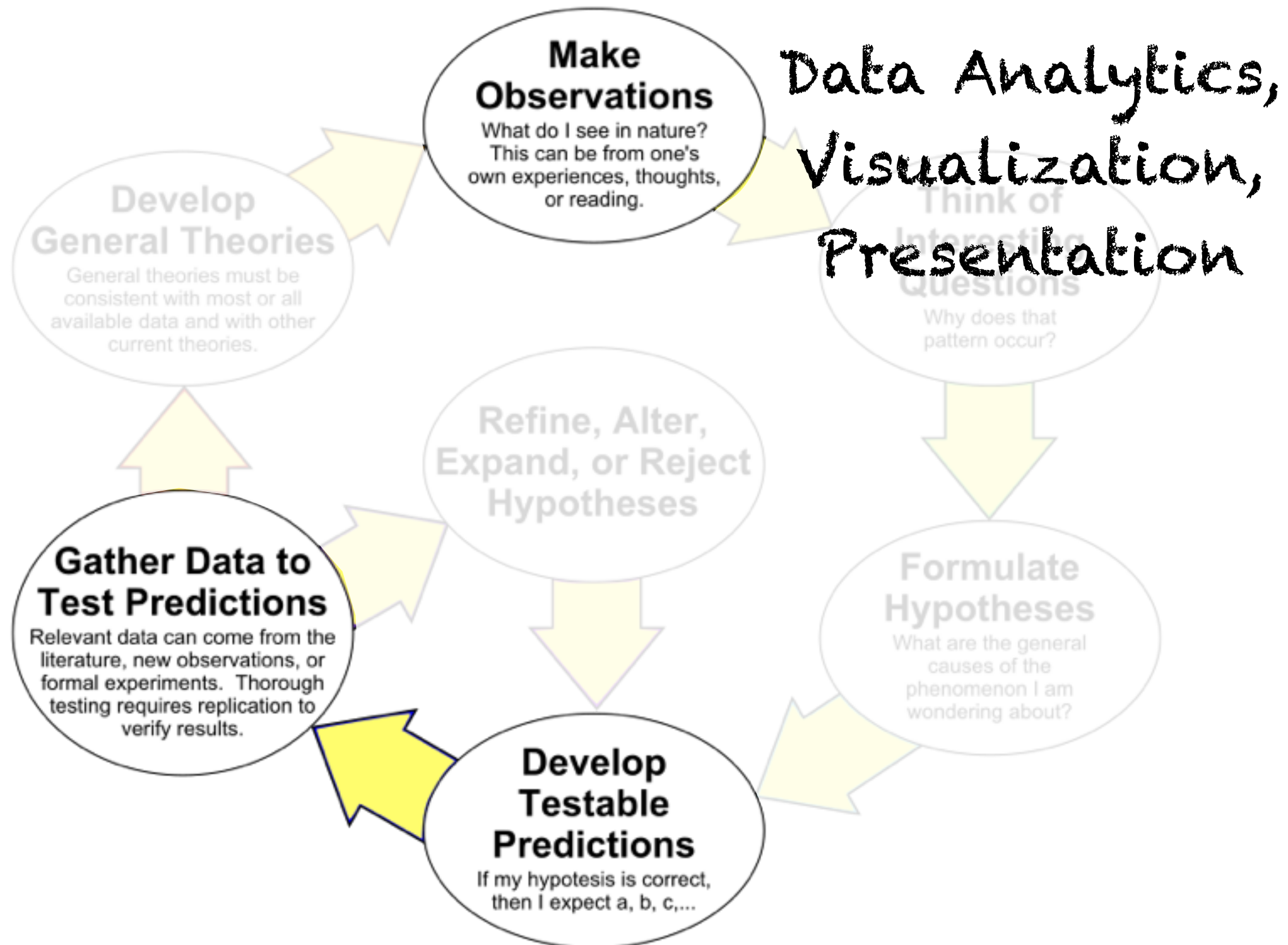


# The Scientific Method

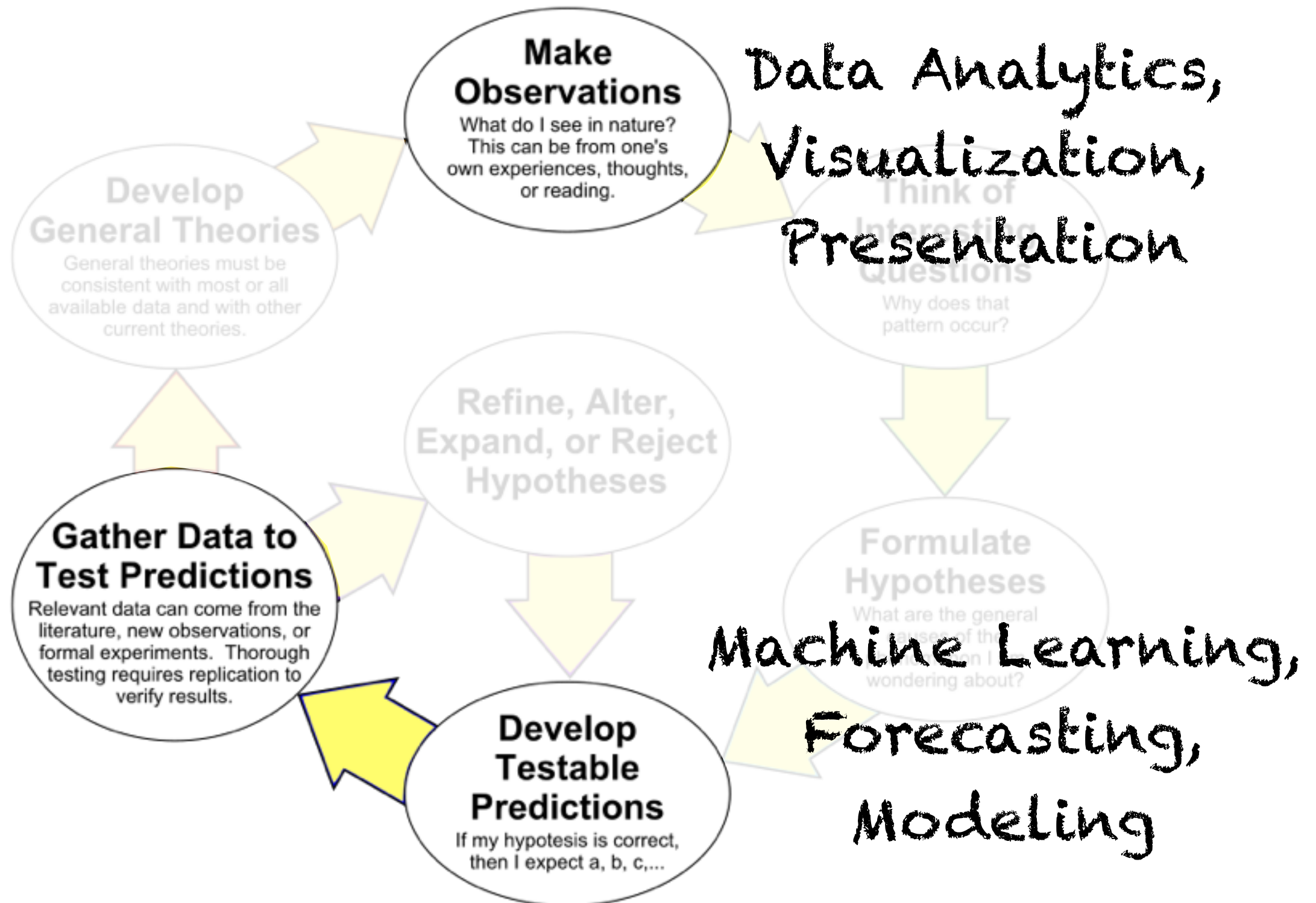




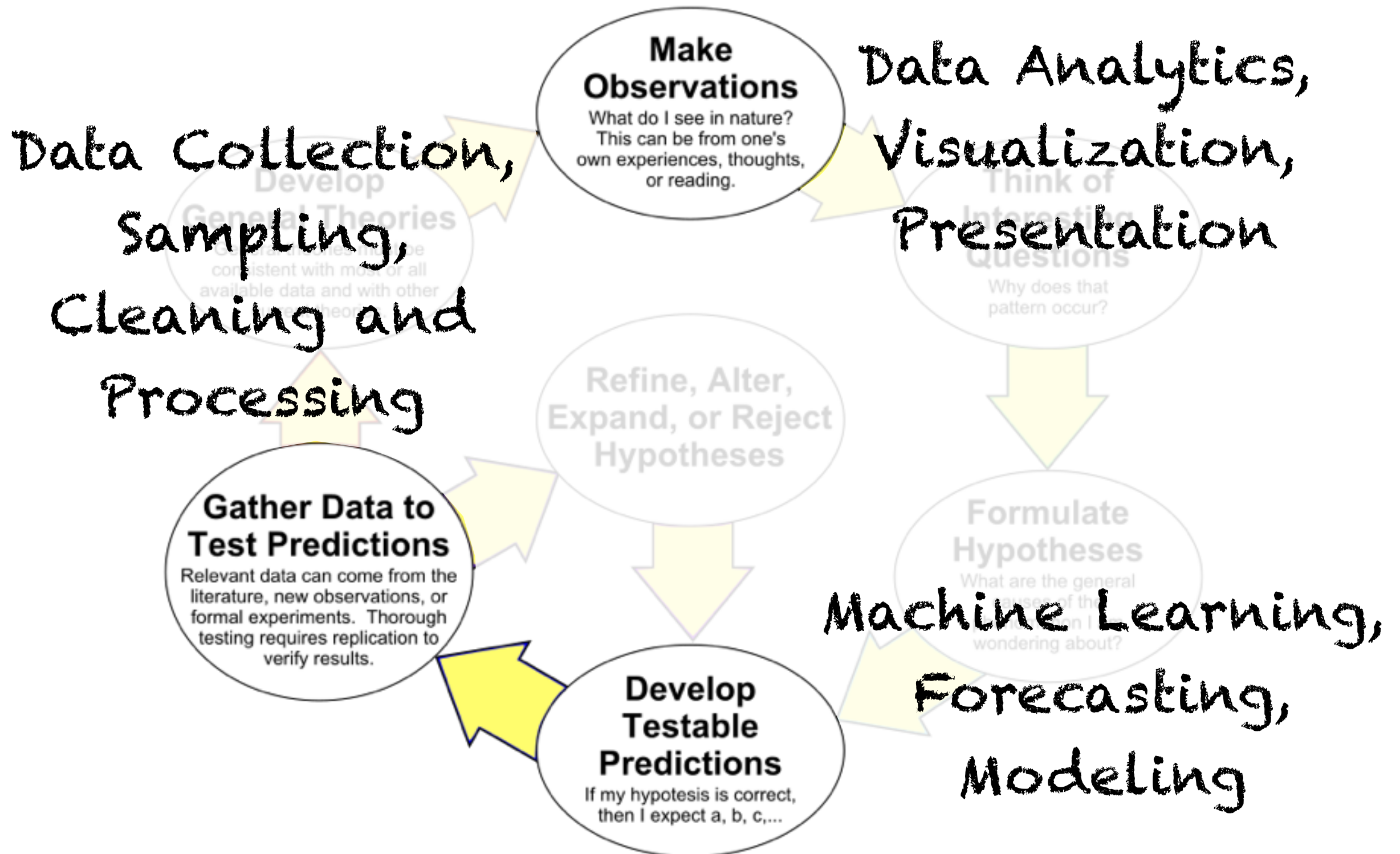
# The Scientific Method



# The Scientific Method

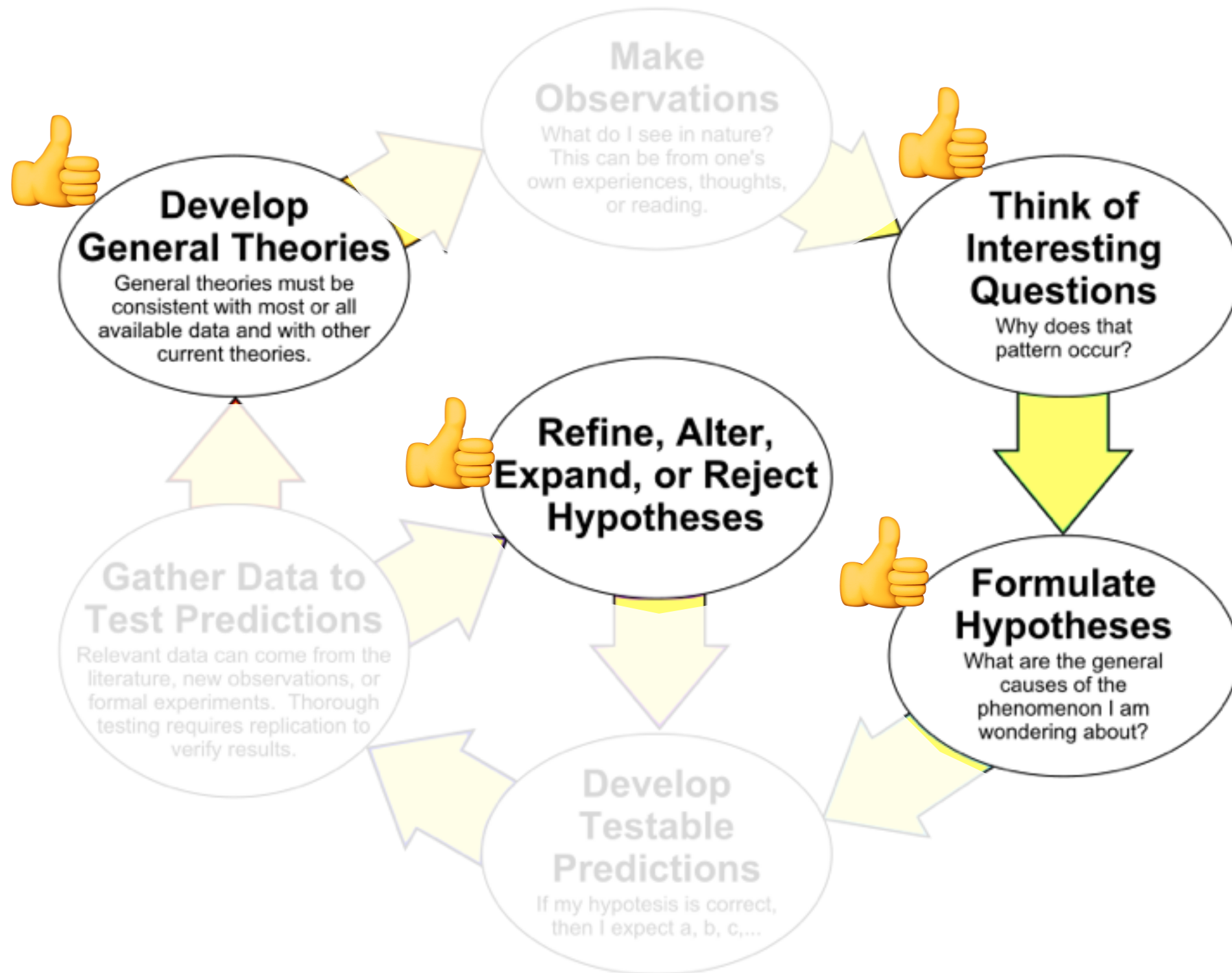


# The Scientific Method

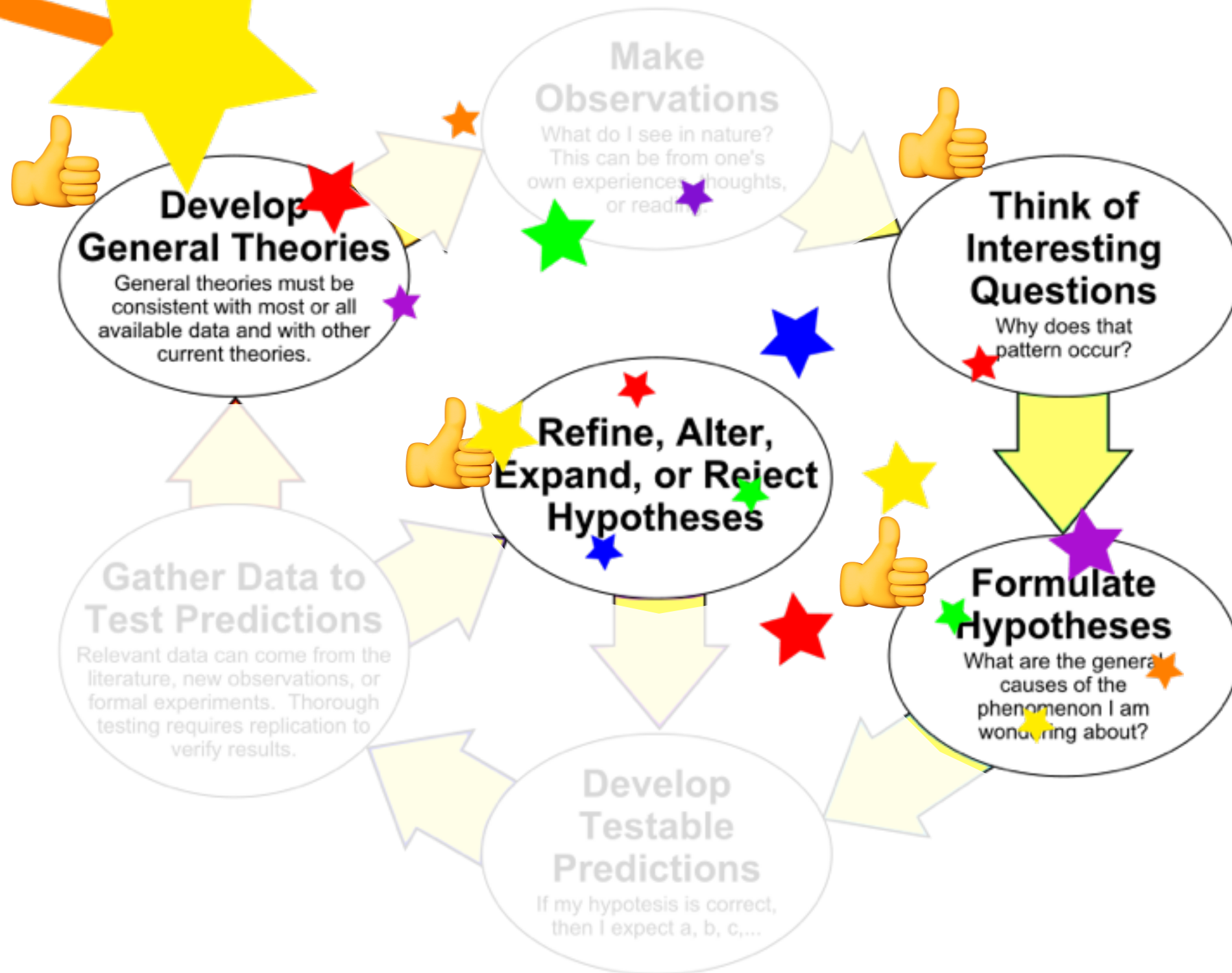




# The Scientific Method



# The Scientific Method



# What is Data Science?





# What is Dads Advice?



**Hang the Christmas lights...**

Data “Science”



# Data “Science”



<https://www.dailydot.com/unclick/state-googled-2017>

<http://nerdgeeks.co/us-state-words-map>



# ata “Science”



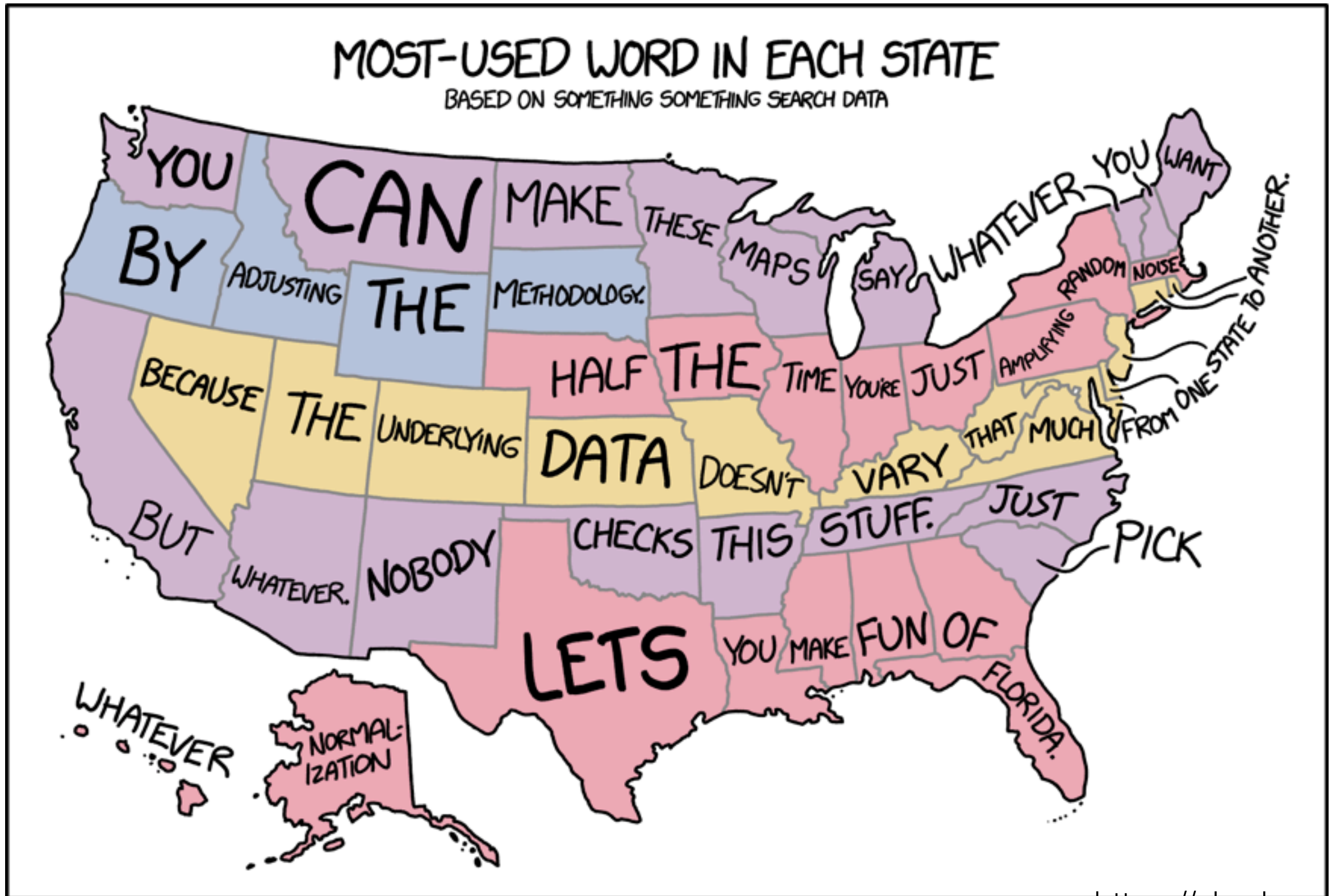
Natalie Delworth



<https://www.dailydot.com/unclick/state-googled-2017>

<http://nerdgeeks.co/us-state-words-map>

# Data “Science”



# Data “Science”

- To be fair...



# Data “Science”

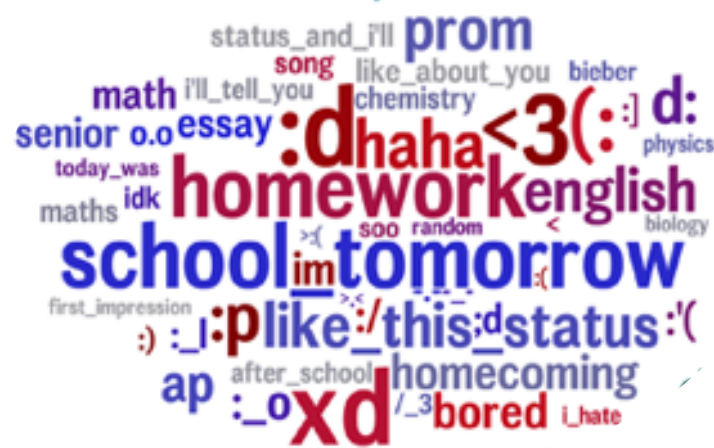
- To be fair...
  - Intuition plays a huge role in the scientific method (“make observations” is Step 1).

# 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”



13-18



19-22

23-29



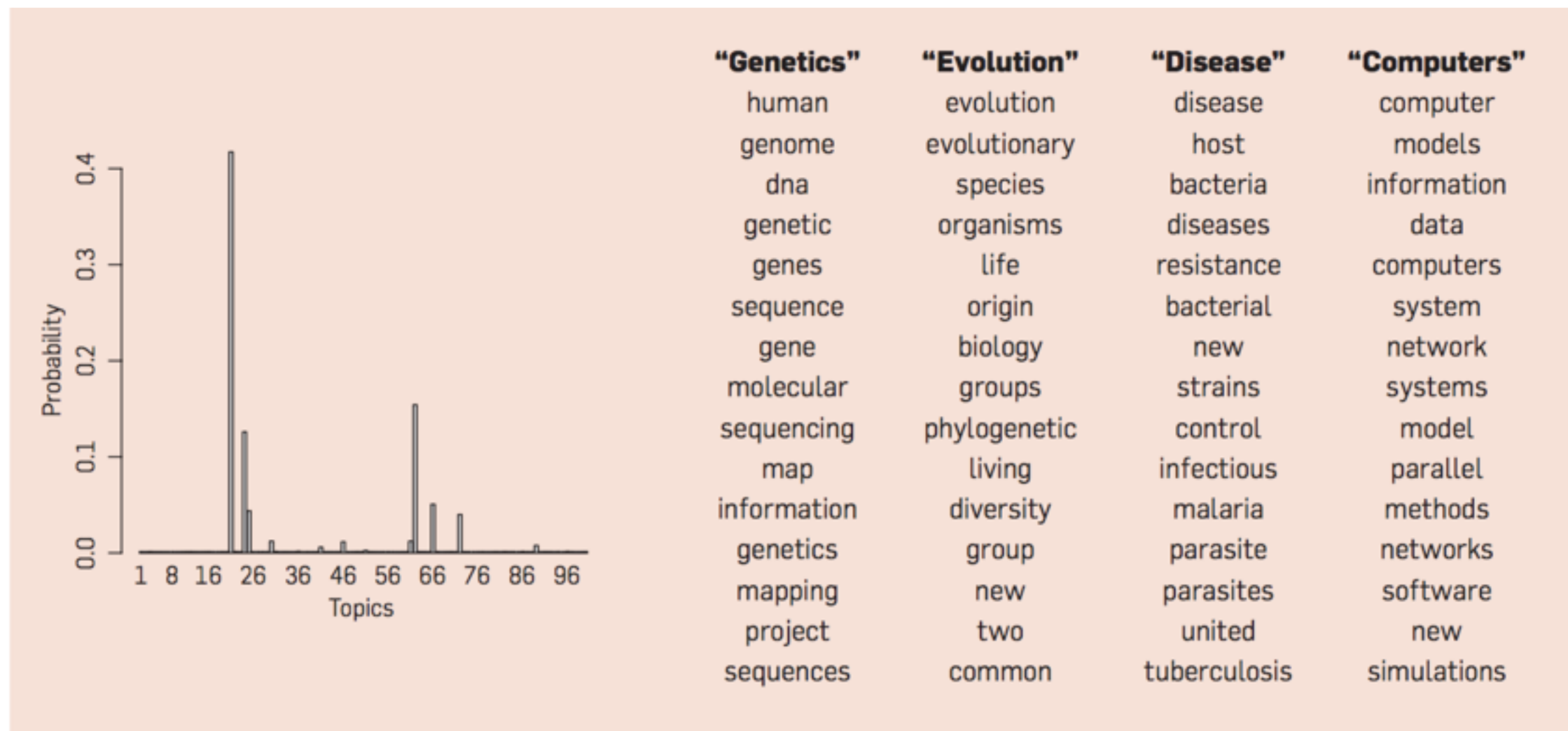
30-65

Facebook posts by age group



# Data “Science”

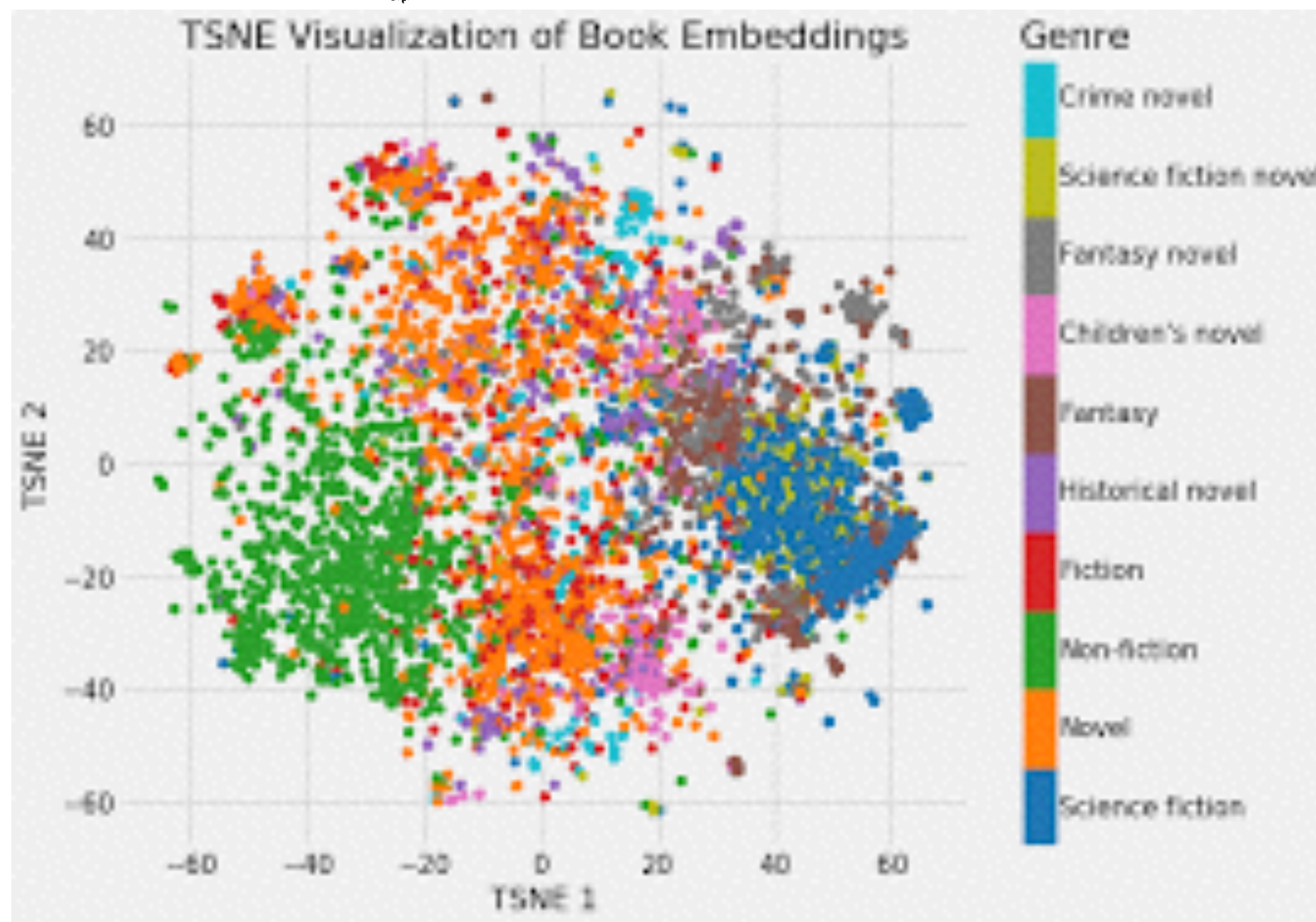
“Eyeballing it”



Frequent topics observed in 17,000 Science articles

# Data “Science”

“Eyeballing it”



# 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
- But!

# Data “Science”

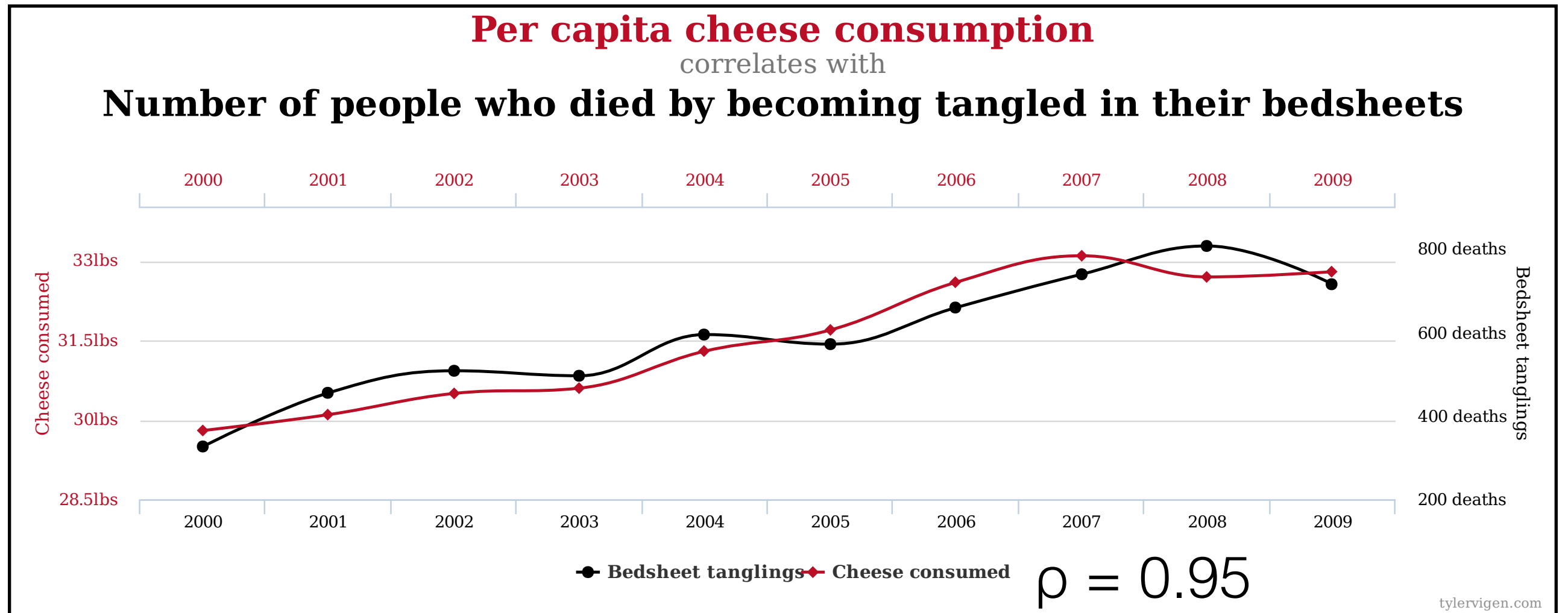
- To be fair...
  - Intuition plays a huge role in the scientific method (“make observations” is Step 1).
  - Exploratory analysis is necessary, its okay to not be all rigor all the time
- But!
  - Exploratory analysis (even when it involves the biggest of data) is meant to \*form\* a hypothesis, not 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, its 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: An argument for multiple comparisons correction**  
Craig M. Bennett<sup>1</sup>, Abigail A. Baird<sup>2</sup>, Michael B. Miller<sup>1</sup>, and George L. Wolford<sup>3</sup>  
<sup>1</sup> Psychology Department, University of California Santa Barbara, Santa Barbara, CA; <sup>2</sup> Department of Psychology, Vassar College, Poughkeepsie, NY; <sup>3</sup> 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<sub>1</sub>-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 included 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

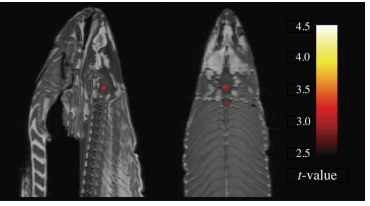
Can we conclude from this data that the salmon is engaging in the perspective-taking task? Certainly not. What we can determine is that random noise in the EPI timeseries may yield spurious results if multiple comparisons are not controlled for. Adaptive methods for controlling the FDR and FWER are excellent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

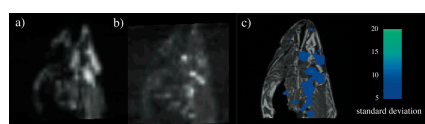


A *t*-contrast was used to test for regions with significant BOLD signal change during the photo condition compared to rest. The parameters for this comparison were  $t(131) > 3.15$ ,  $p(\text{uncorrected}) < 0.001$ , 3 voxel extent threshold.

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<sup>3</sup> 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 ( $p = 0.25$ ).

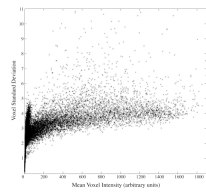
### VOXELWISE VARIABILITY



To examine the spatial configuration of false positives we completed a variability analysis of the fMRI timeseries. On a voxel-by-voxel basis we calculated the standard deviation of signal values across all 140 volumes.


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<sub>1</sub>-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.

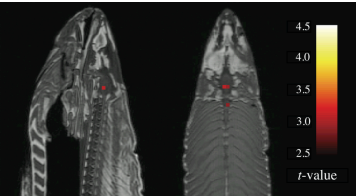


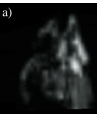
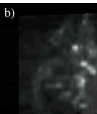
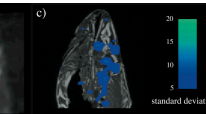
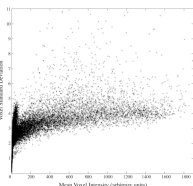
Neural correlates of interspecies perspective taking in the post-mortem Atlantic Salmon

# Data “Science”

 **Neural correlates of interspecies perspective taking in the post-mortem Atlantic Salmon: An argument for multiple comparisons correction**  
Craig M. Bennett<sup>1</sup>, Abigail A. Baird<sup>2</sup>, Michael B. Miller<sup>1</sup>, and George L. Wolford<sup>3</sup>  
<sup>1</sup> Psychology Department, University of California Santa Barbara, Santa Barbara, CA; <sup>2</sup> Department of Psychology, Vassar College, Poughkeepsie, NY; <sup>3</sup> 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

**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<sup>3</sup> 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 ( $p = 0.25$ ).

**VOXELWISE VARIABILITY**  
a)  b)  c)   
To examine the spatial configuration of false positives we completed a variability analysis of the fMRI timeseries. On a voxel-by-voxel basis we calculated the standard deviation of signal values across all 140 volumes.  
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<sub>1</sub>-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.  


**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.

**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.

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



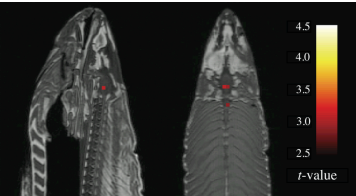
# Data “Science”

 **Neural correlates of interspecies perspective taking in the post-mortem Atlantic Salmon: An argument for multiple comparisons correction**  
Craig M. Bennett<sup>1</sup>, Abigail A. Baird<sup>2</sup>, Michael B. Miller<sup>1</sup>, and George L. Wolford<sup>3</sup>  
<sup>1</sup> Psychology Department, University of California Santa Barbara, Santa Barbara, CA; <sup>2</sup> Department of Psychology, Vassar College, Poughkeepsie, NY; <sup>3</sup> 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.

**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<sup>3</sup> 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 ( $p = 0.25$ ).

**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<sub>1</sub>-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 included 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**

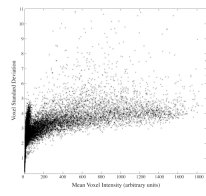
Benjamini Y and Hochberg Y (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society: Series B*, 57:289-300.

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

**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<sub>1</sub>-weighted image.

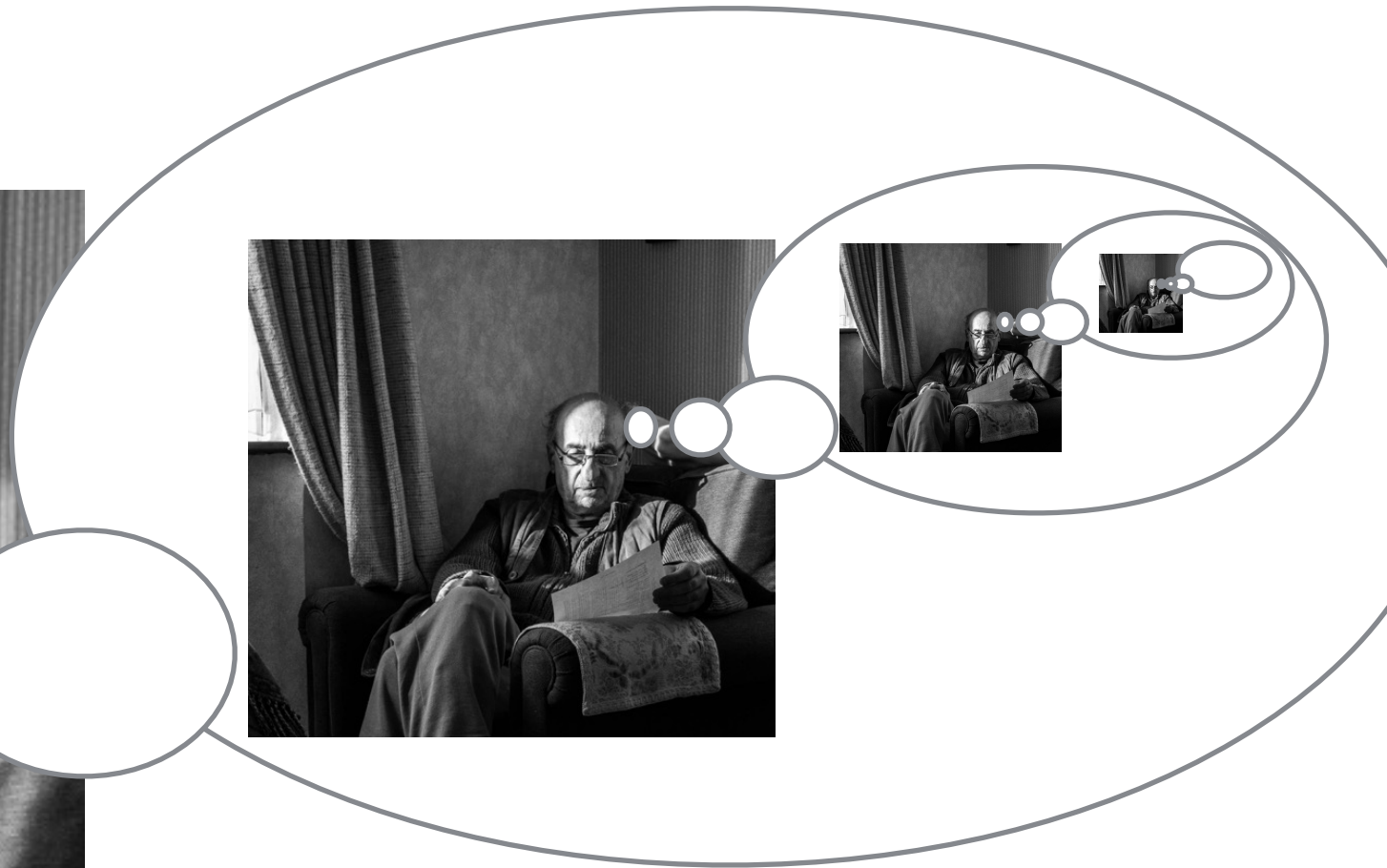
To investigate this effect in greater detail we conducted a Pearson correlation to examine the relationship between the signal in a voxel and its variability. There was a significant positive correlation between the mean voxel value and its variability over time ( $r = 0.54$ ,  $p < 0.001$ ). A scatterplot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





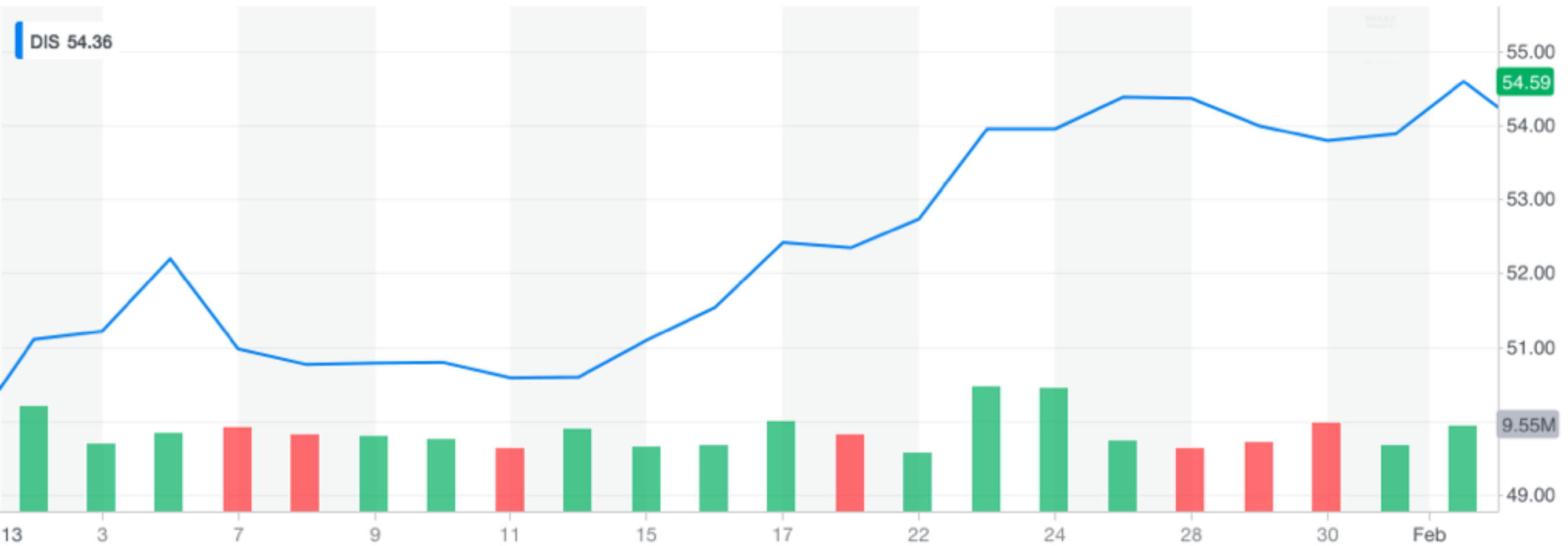


Roses are red,  
Violets are blue.





# “Data” Science



# “Data” Science

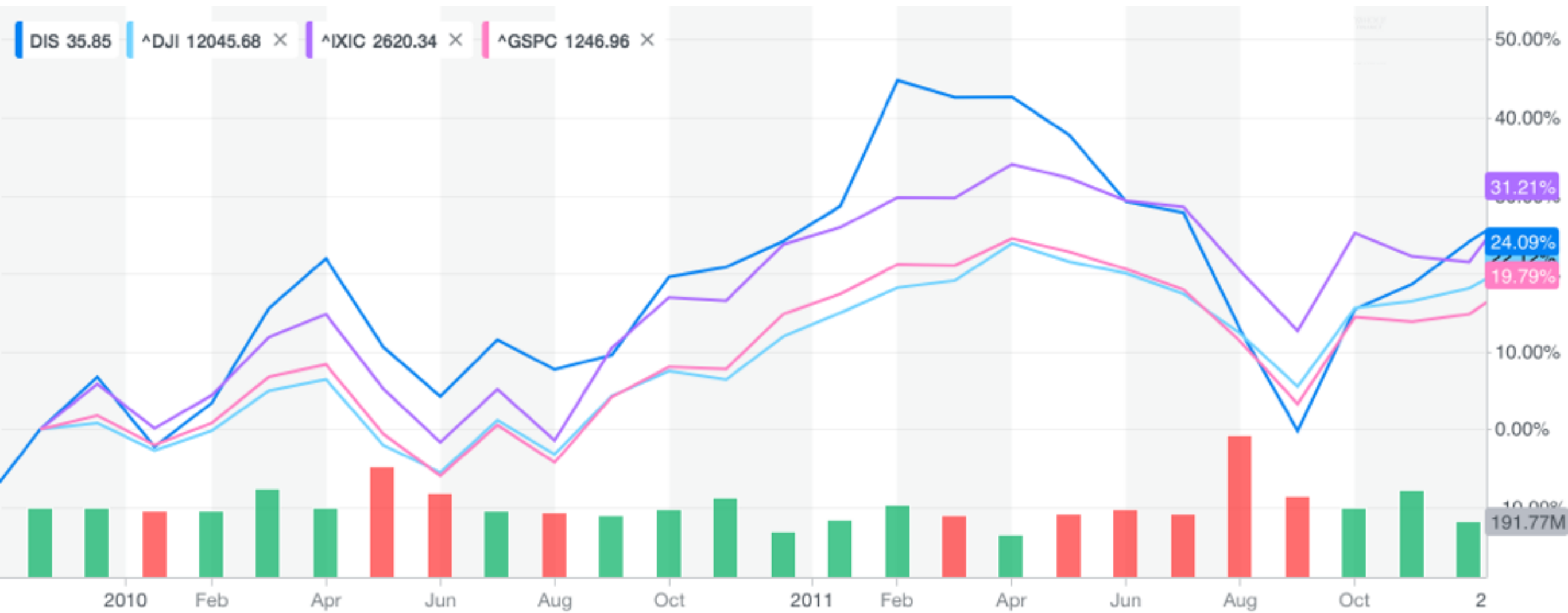


# “Data” Science



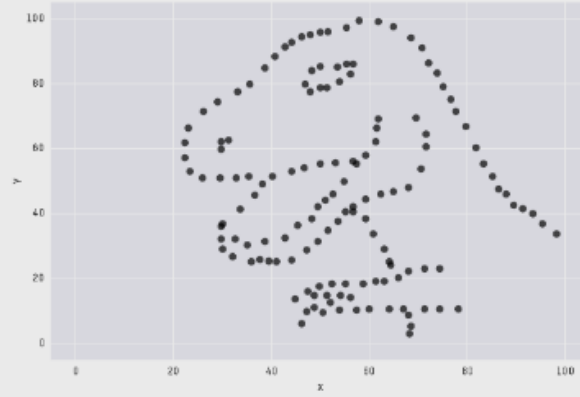


# “Data” Science

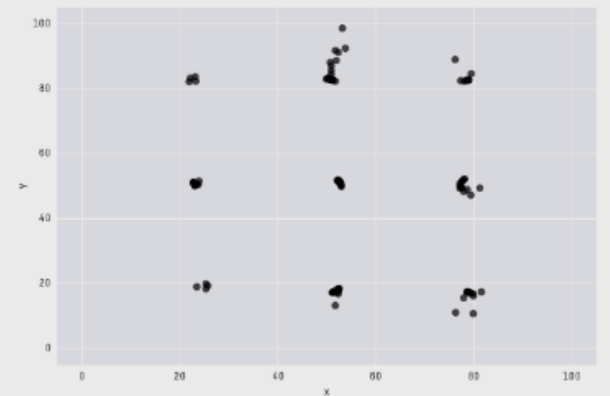
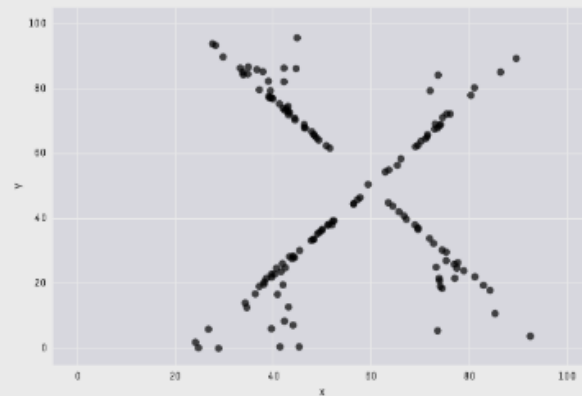
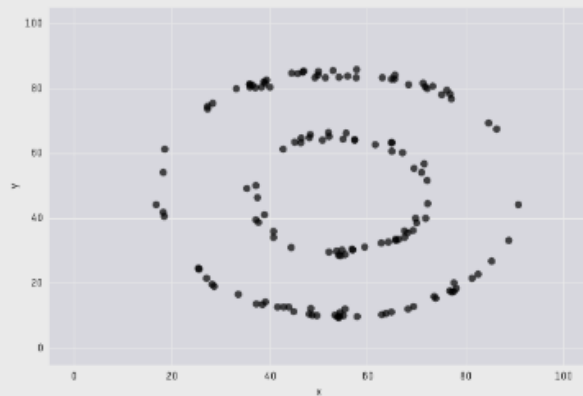
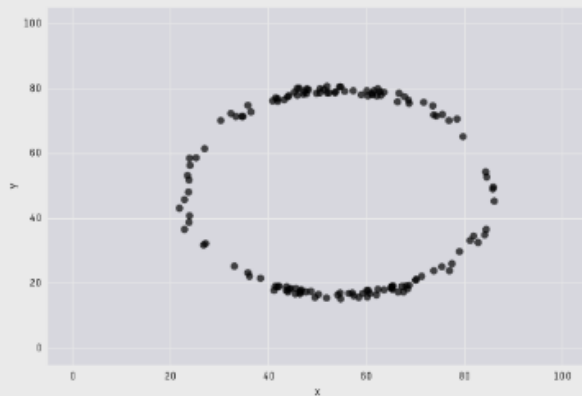
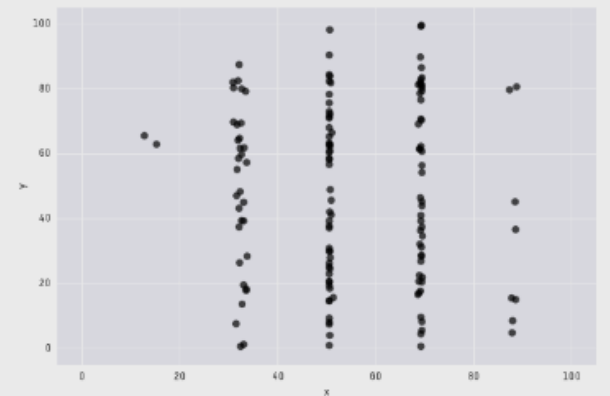
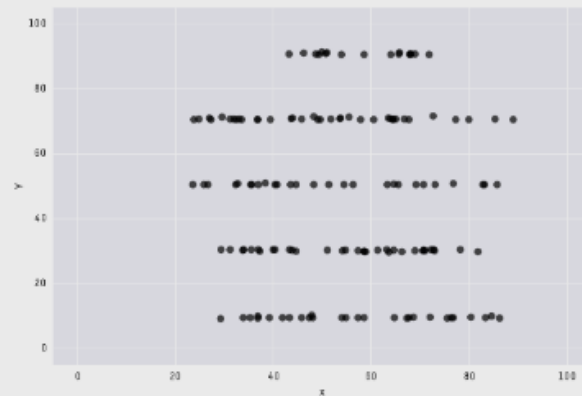
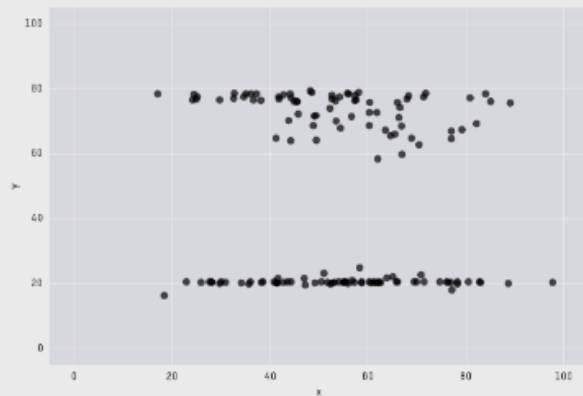
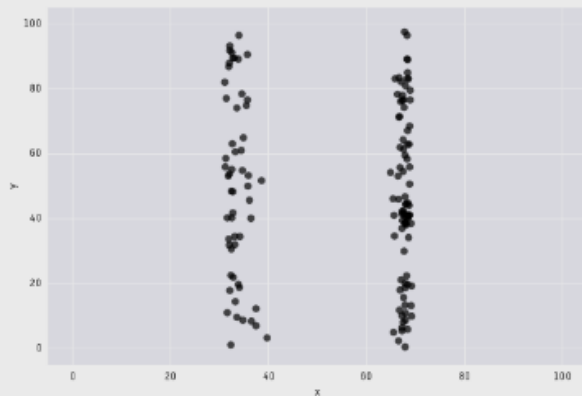
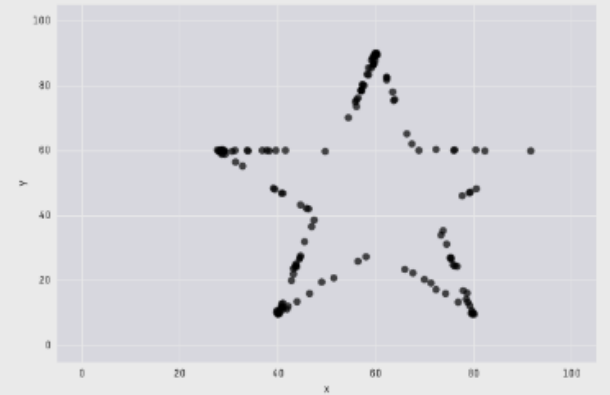
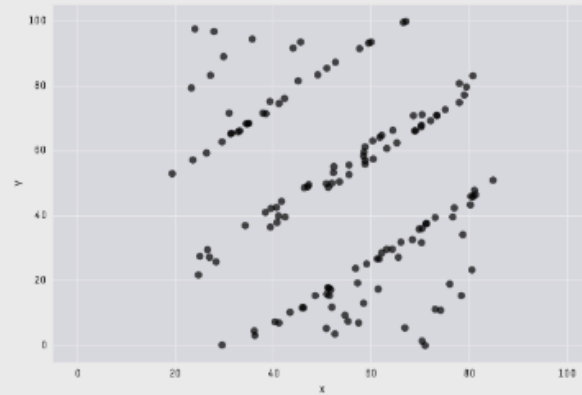
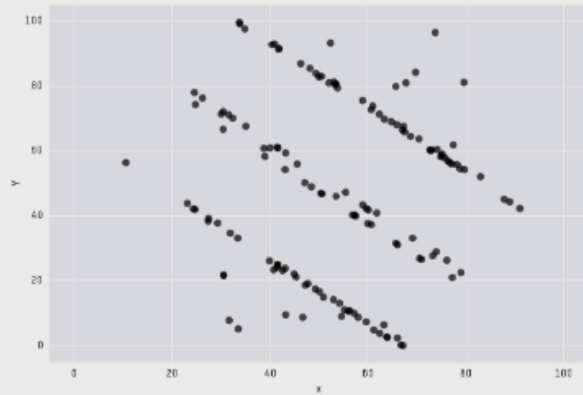
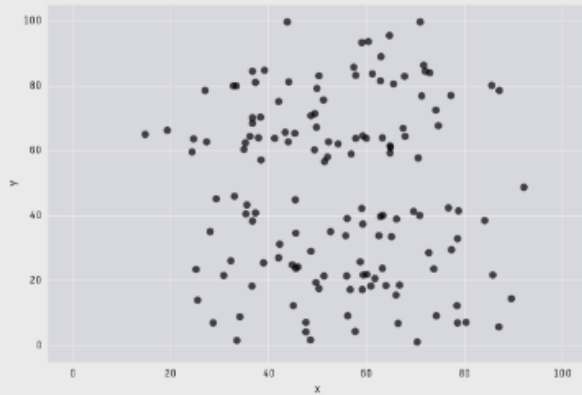


# “Data” Science

```
X Mean: 54.26  
Y Mean: 47.83  
X SD   : 16.76  
Y SD   : 26.93  
Corr.  : -0.06
```



X Mean: 54.26  
 Y Mean: 47.83  
 X SD : 16.76  
 Y SD : 26.93  
 Corr. : -0.06



<https://blog.revolutionanalytics.com/2017/05/the-datasaurus-dozen.html>

shout out **Kevin Jin** for sharing this last year! :)

# “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.

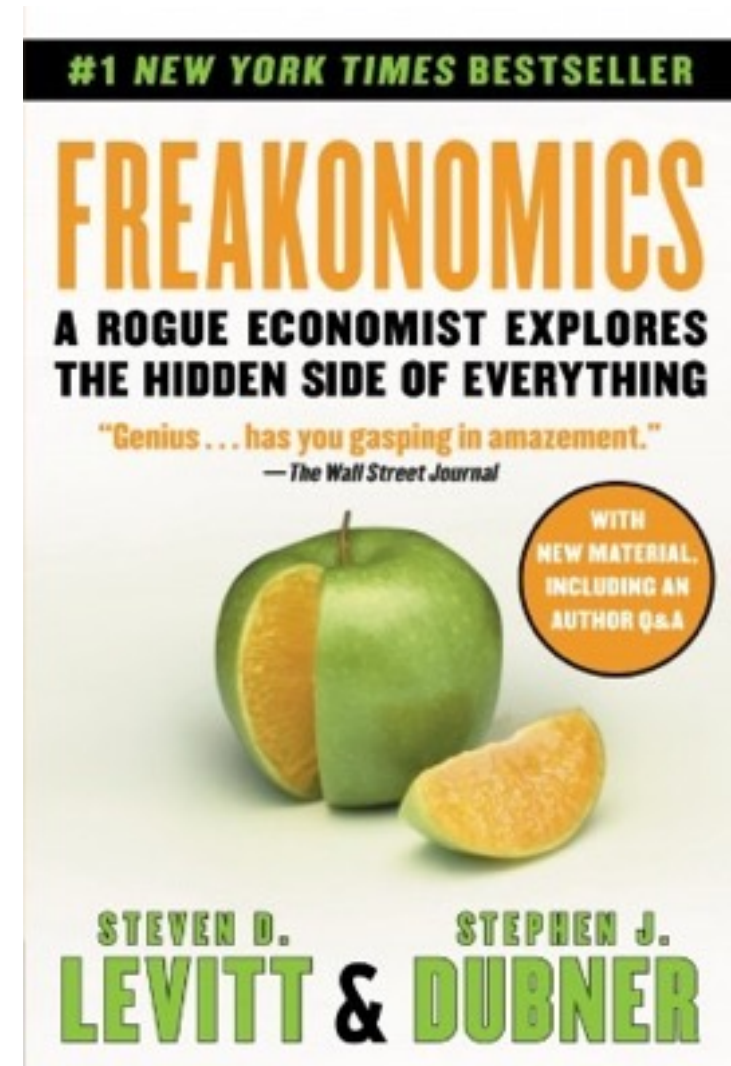


# “Data” Science

- 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



# “Data” Science

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

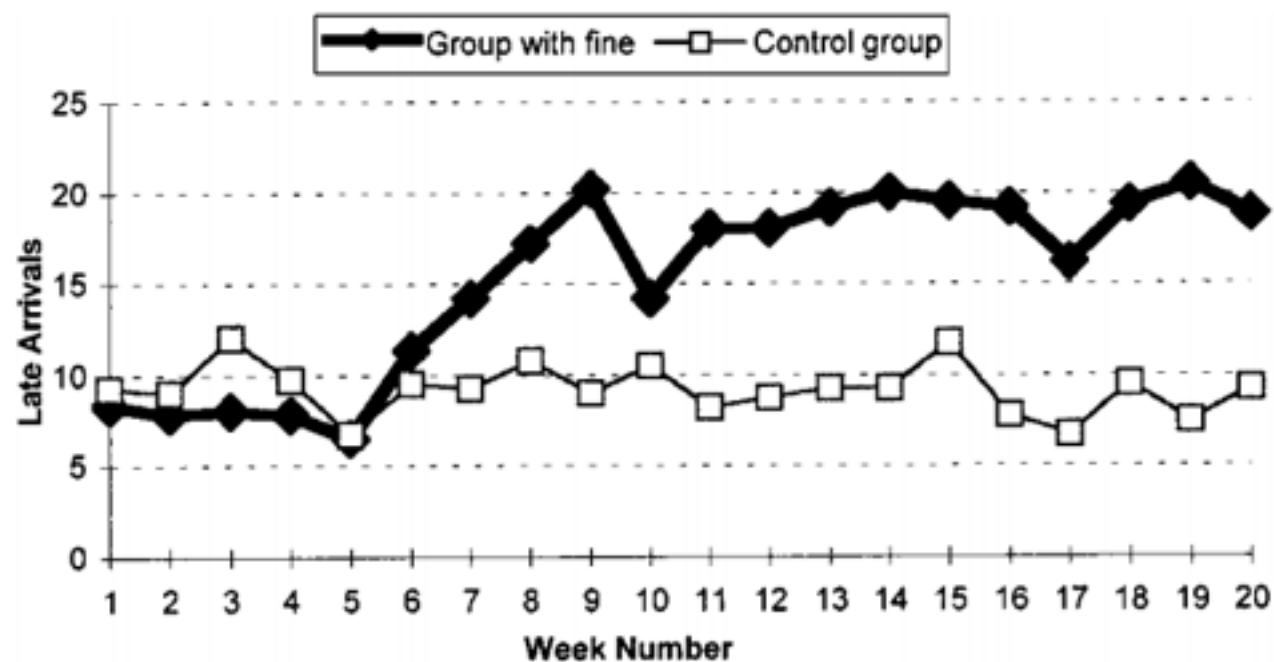
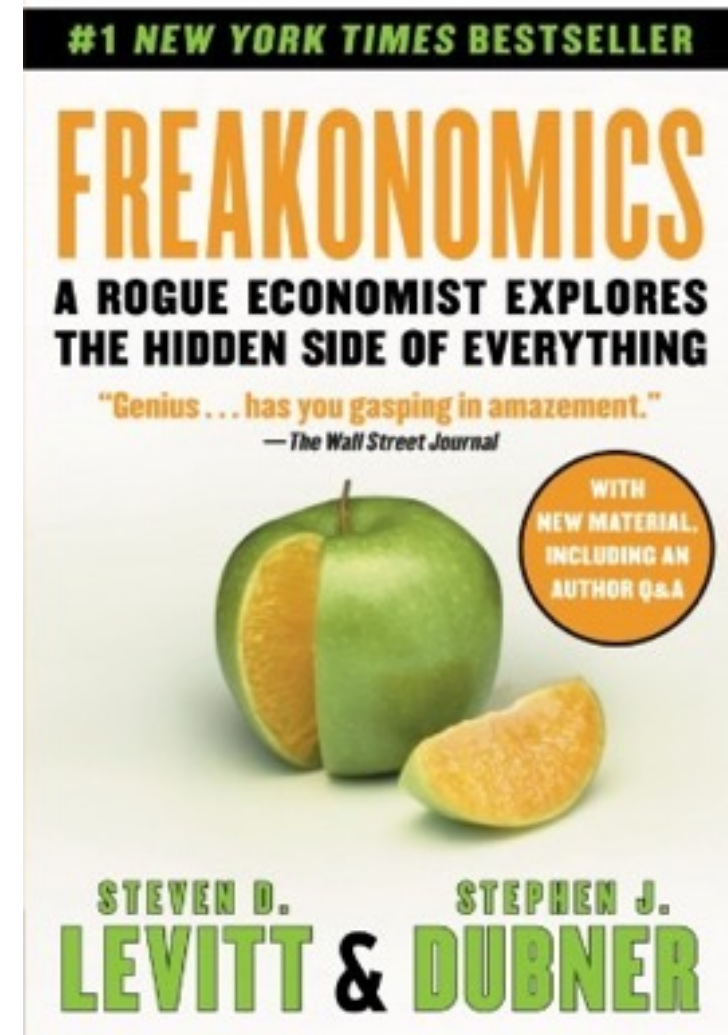
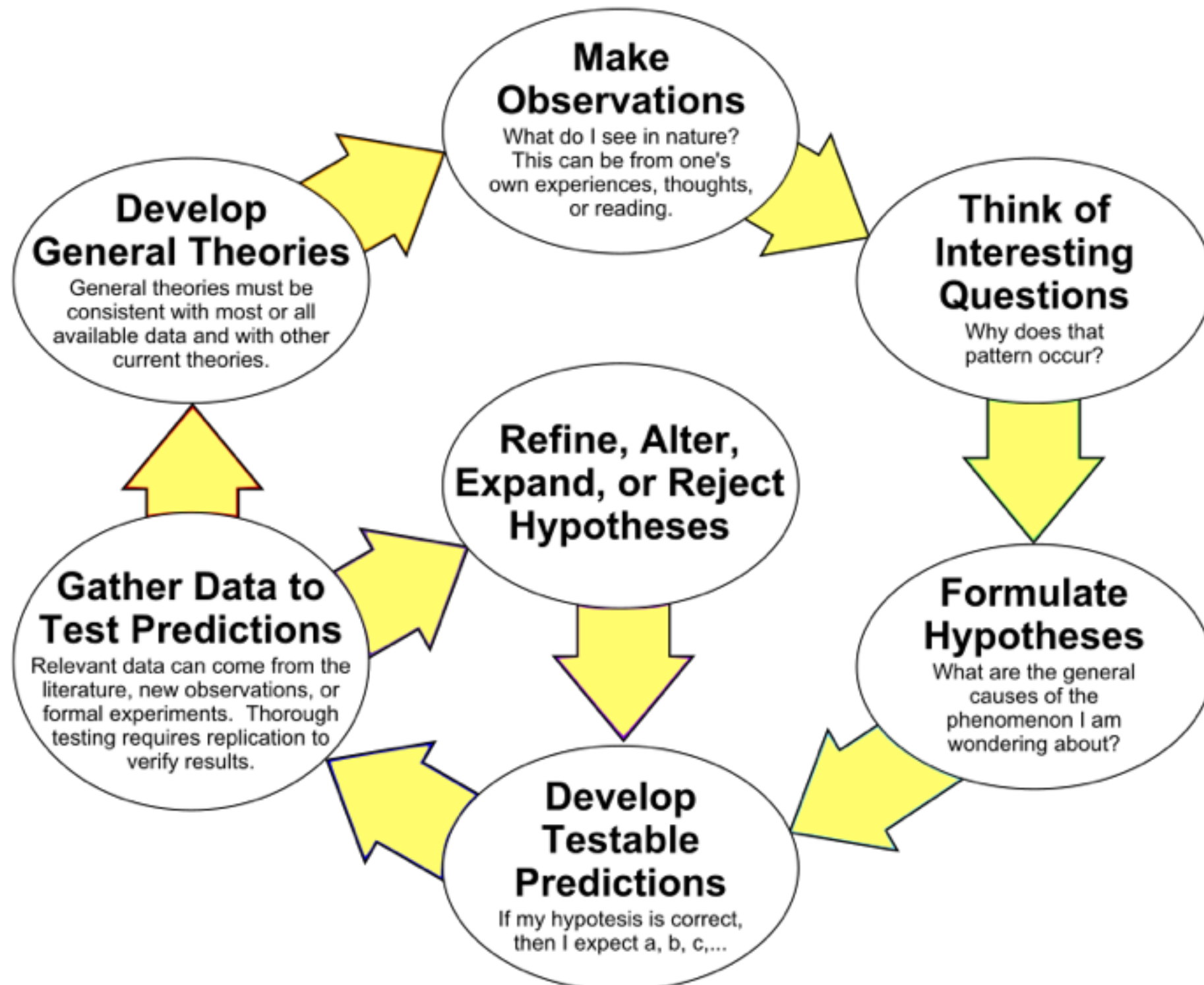


FIGURE 1.—Average number of late-coming parents, per week



# Data! Science!





CSCI 1951A

What is ~~Data Science~~?

## January

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	

## 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	29

## 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	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
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	

## 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	29

## 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	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		

This. Right Here, Right Now.

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

- Databases for Data Scientists: Entity-Relationship (ER) Diagrams, SQL **[Assignment 1]**
- Web Crawling, API Calls **[Assignment 2]**
- Data Cleaning and Normalization
- Crowdsourcing
- Working at Scale: MapReduce, Google Cloud **[Assignment 3]**



## January

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	

## 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	29

## 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	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		

- Probability and Statistics
- Hypothesis Testing **[Assignment 4]**
- P-Values (and their pitfalls)
- T-Tests, Chi-Squared Tests, Regression
- Working with stats\_models

## January

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	

## 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	29

## 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	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		

- Intro ML: feature representations, loss functions
- Types of models: supervised vs. unsupervised learning
- Clustering with K-Means **[Assignment 5]**
- Regression revisited, prediction vs. hypothesis testing
- Overfitting and regularization
- Working with sklearn

## January

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	

## 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	29

## 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	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 Visualization in D3 **[Assignment 6]**
- Just enough html and javascript to do D3 :)

## January

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	

## 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	29

## 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	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		

- Natural Language Processing 101 **[Assignment 7]**
- ML Fairness
- Matrix Factorization and Recommender Systems
- Deep Learning 101

## January

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	

## 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	29

## 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	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		

- Feb 6: Project Proposals
- Feb 27: Data: Done! (Scraped, Cleaned, Databased). No changing plans after this.
- March 19: Stats Deliverable. Initial analysis...i.e. evidence your first idea was wrong/won't work. ;)
- April 2: Mid-Semester Feedback
- April 9: Viz Deliverable...i.e. when you realize something about your data you probably should have known already
- May 7: Final Project Due. Poster Day



# Grading

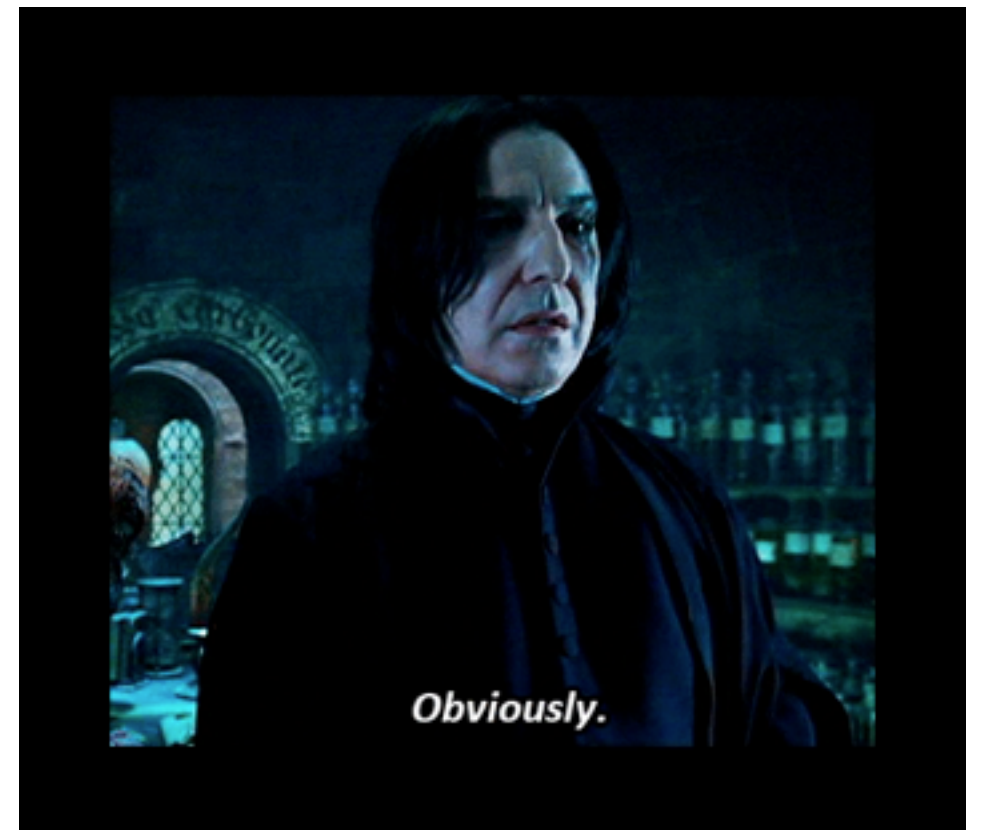
- 50% Assignments (~7% each)
- 30% Final Project
- 10% Labs
- 10% Attendance/Clickers (must attend 2/3 of classes)

# Late Days

- Assignments are due at 11:59 pm on the listed due date
- 7 late days total; no maximum per assignment
- 20% penalty for each additional day late
- No late days for Final Project deliverables (incl. intermediate deliverables)
- Deans Notes/SEAS? -> talk to Ellie
- Any other extension requests? -> No.

# Collaboration

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



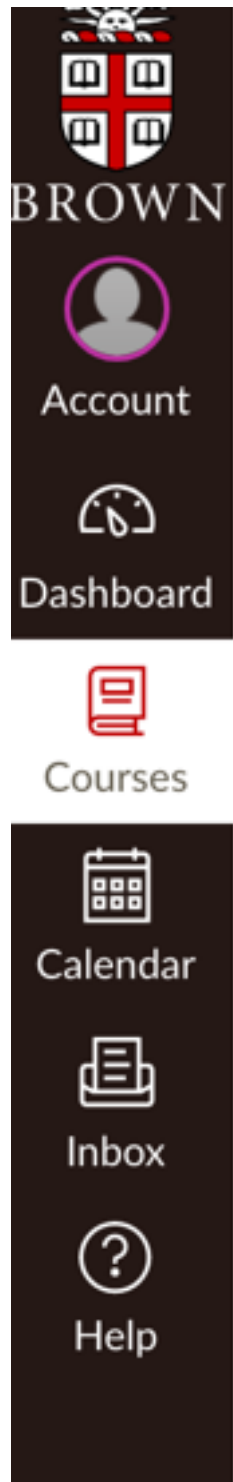
To Do Now



# 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

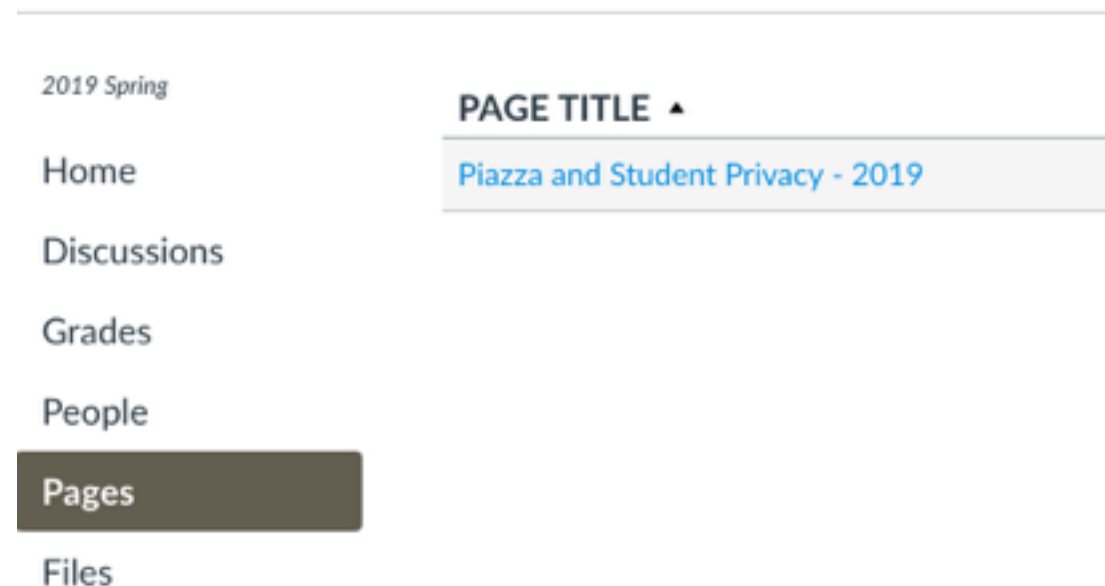
Chat

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.



# To Do Now

- Hours are starting Sunday! Go say hi to your staff...
- SQL assignment will be released tomorrow

# To Do Now

- Start brainstorming final projects and forming groups! Project group mixer soon, TBD.
- Things to consider:
  - do we want to do the same thing? (duh)
  - capstone
  - do we work at the same pace?
  - do we work during the same hours?
  - do we communicate the same way?
  - do I even like this person...?



Thank you!  
Questions?