# About us!

also my collaborators…



**Jenny Smetzer**



**Chester Ismay**

# About you!

Say hi to your nearest neighbors!
You'll be learning together!

# My Context for moderndive

**My students:**
- Undergraduate-only liberal arts college
- Service intro stats course for all majors, all years
- Calculus is a pre-req only in name
- 13 weeks x (3 x 70min lectures + 75min lab)
- 29/40 had never coded in R prior

**My goals:**
- Goal 1: Sampling for inference
- Goal 2: Modeling with regression

# Getting from Point A to Point B

## via the `tidyverse`

**Point A:**
Modal 1st time
stats student

**Point B:**
Two goals

😐

1. Sampling for inference
2. Modeling with regression

Calculus?
😬 thru 🤮

Coding?
😱 & 🤔

# What is the tidyverse?



Core: Loaded with library(tidyverse)

Commonly-Used: Installed with tidyverse, but not loaded by default with library(tidyverse)

- ggplot2 for data visualization
- dplyr for data wrangling
- readr for data importing

# Why tidyverse? Some principles

From [tidy tools manifesto]():     Say what?

1. Reuse existing data structures
2. Compose simple functions with the pipe
3. Embrace functional programming
4. Design for humans

1. Don't reinvent the wheel!
2. Break down tasks step-by-step!
3. What is the [goal]() of your code?
4. Make code understandable

# Why tidyverse for stats newbies?

- *In my opinion* it's easier to learn than base R. [Others too](#).
- It scales. You leverage an entire ecosystem of online developers and support: Google & StackOverflow
- Satisfy learning goals *while learning tools they can use beyond the classroom*

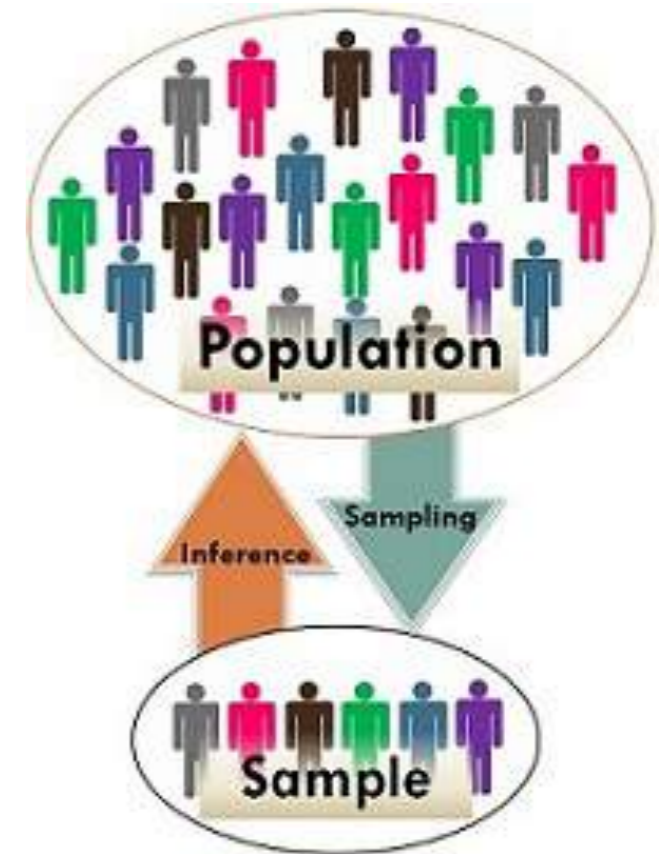# Goal 1: Sampling for Inference

1. Tactile Sampling → 2. Virtual Sampling → 3. Theoretical

**Population**



```
Console ~/
> library(moderndive)
> bowl
# A tibble: 2,400 × 2
   ball_ID color
     <int> <chr>
1        1 white
2        2 white
3        3 white
4        4 red
5        5 white
6        6 white
7        7 red
8        8 white
9        9 red
10      10 white
# … with 2,390 more rows
>
```
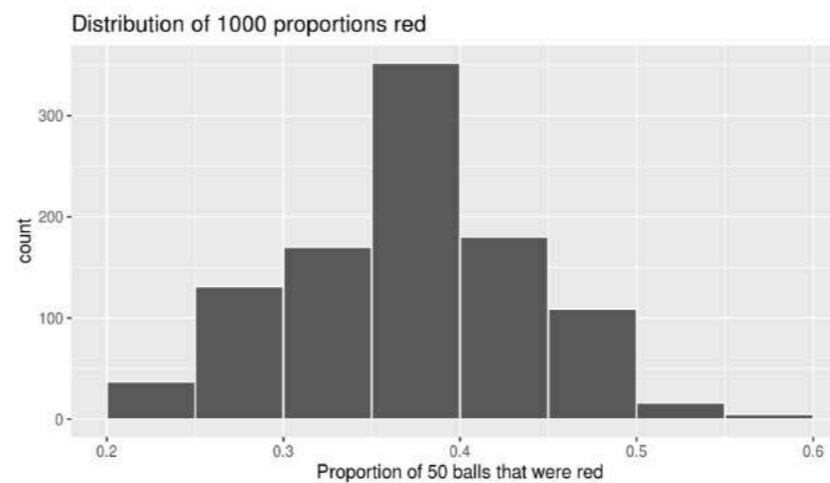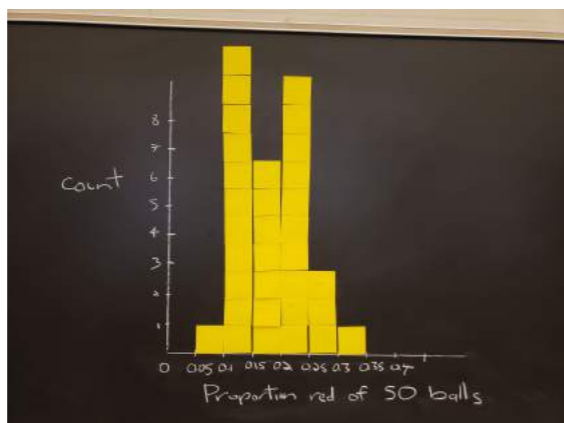


**Sample**



```
Console ~/
> bowl %>%
+   rep_sample_n(size = 50, reps = 1)
# A tibble: 50 × 3
# Groups:   replicate [1]
   replicate ball_ID color
       <int>   <int> <chr>
1          1     226 white
2          1    1304 red
3          1    1230 white
4          1     984 white
5          1      68 white
6          1    1965 white
7          1     431 white
8          1    1184 white
9          1    1610 red
10         1     978 white
# … with 40 more rows
>
```
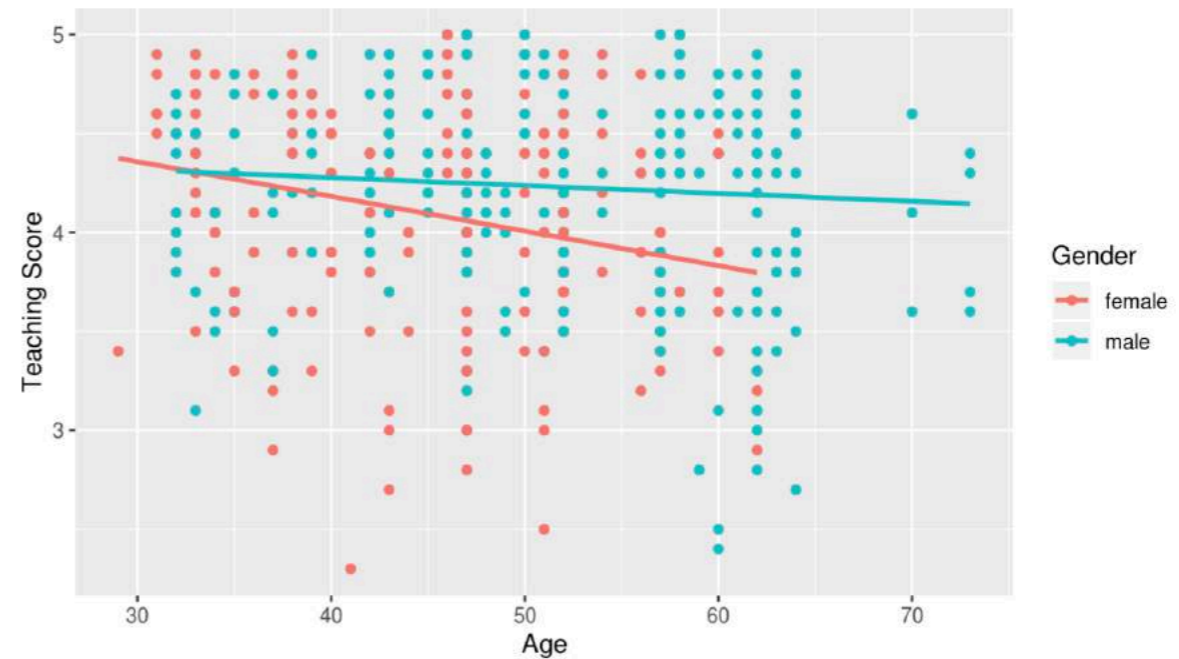
**Sampling Distributions & Standard Errors**



Distribution of 1000 proportions red



$$SE = \sqrt{\frac{p(1-p)}{n}}$$

9

# Goal 2: Modeling with Regression

## 1. Data



## 2. Exploratory Data Analysis



## 3. Regression Coeff



Early: Descriptive regression

**Sampling**

**Conf Int's & Hyp Tests**

## 4. Regression Table



Later: Inference for Regression

# End Deliverable

Final project that "plays the whole game"
of data/science pipeline:



Example template given to students this semester,
based on work by Alexis Cohen, Andrianne Dao, & Isabel
Gomez last semester.

# Schedule

See Google Doc available at bit.ly/USCOTS2019

# Keep in mind throughout…

- You are not currently you, but you are currently your students *as best as you can imagine*.
- In other words, these exercises are meant for your students!
- Ultimately where do I start? *Start small!*

# Questions?                    Let's Go!

# Activity: Your Birthday

# For Every Chapter…

Slides on:
1. What are we doing ❓
2. Why are we doing this 🤔
3. Opinions
4. Potential pitfalls ⚠️

Followed by activities:
1. Chalk talk, pen/paper, or tactile exercise
2. Replicating exercise on computer
3. Exercise
4. Discussion

# Chapter 2: Exploring data

1. What are we doing **?**
   - Getting used to workspace via data exploration
2. Why are we doing this 🤔
   - Getting them over initial 😱 of coding
3. Our opinions
   - Stress importance of looking at RAW data values. Removing these layers of abstraction.
4. Potential pitfalls ⚠️
   - Difference of R vs RStudio
   - Installing/loading packages
   - Error messages, warning messages, regular messages
   - Coding: both student self doubt & lowered instructor expectations

# Chapter 3: Visualizing Data

1.  What are we doing ❓
    - Creating (colored) scatterplots, histograms, boxplots
2.  Why are we doing this 🤔
    - Equip students with tools for both our goals
    - [Exploratory data analysis!!!](#)
3.  Our opinions
    - Viewing all graphics through lens of the Grammar of Graphics (via `ggplot2`)
4.  Potential pitfalls ⚠️
    - Histograms & boxplots involve transformations of raw values
    - Coding ramps up: Reassure students! Encourage them to not code from scratch, rather copy/paste/tweak

# Chapter 4: Data Wrangling

1. What are we doing **?**
   - Learning the pipe operator **%>%**
   - Wrangling/transforming data
2. Why are we doing this 🤔
   - Equip students with tools for both our goals
3. Our opinions
   - To *completely* shield students from *any* data wrangling is to betray [true nature of work in our fields]
4. Potential pitfalls ⚠️
   - How much wrangling should you require vs you curate yourself?

# Chapter 5: "Tidy" data

1. What are we doing **?**
   - tidyverse gets its name from fact that all inputs/outputs are assumed to be *tidy data frames*
   - Importing data via `readr::read_csv()`
2. Why are we doing this 🤔
   - Students have their own Excel/Google Sheets data
   - Will have to convert from wide to tidy/long format
3. Our opinions
   - This chapter can be skipped if
     A. You only provide tidy/long data
     B. You have your students [publish .csv](#) files to Google Sheets
4. Potential pitfalls ⚠️
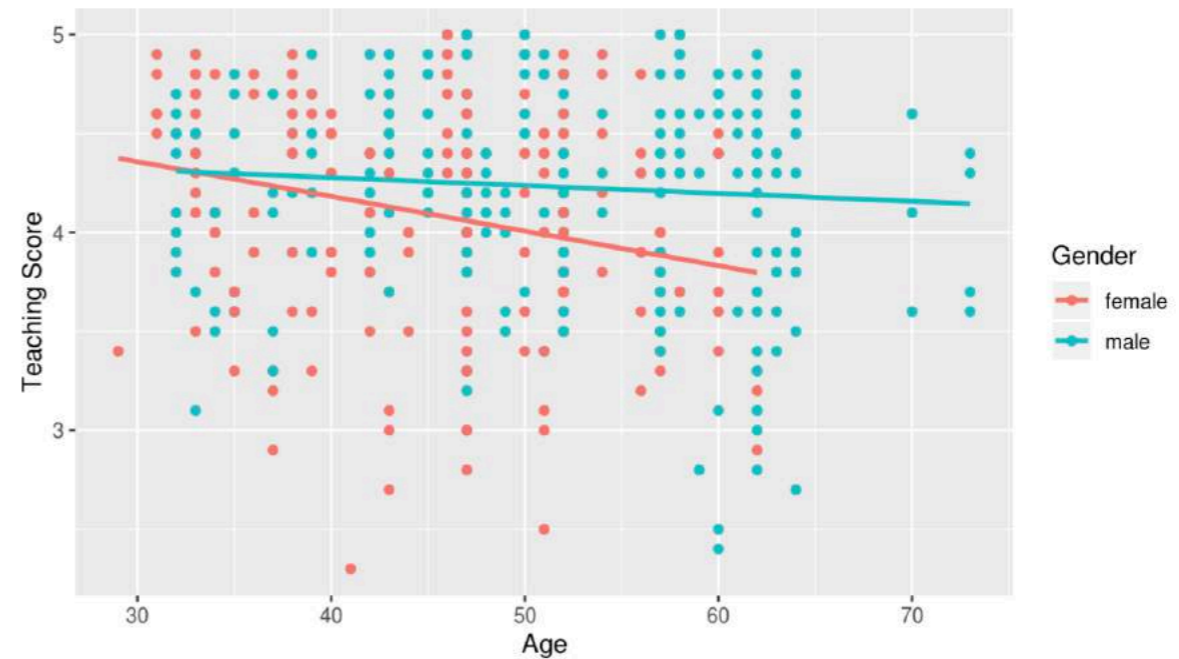   - Working directories!

# Chapter 6: Simple regression

# Goal 2: Modeling with Regression
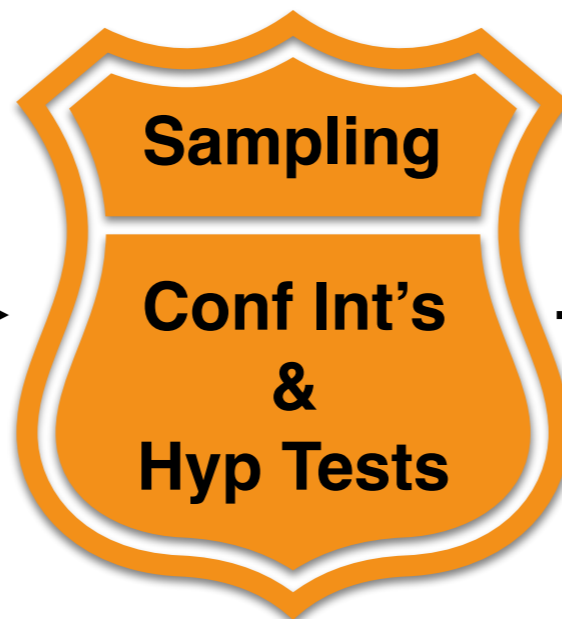
## 1. Data

## 2. Exploratory Data Analysis



## 3. Regression Coeff



Early: Descriptive regression

## 4. Regression Table



Sampling

Conf Int's
&
Hyp Tests
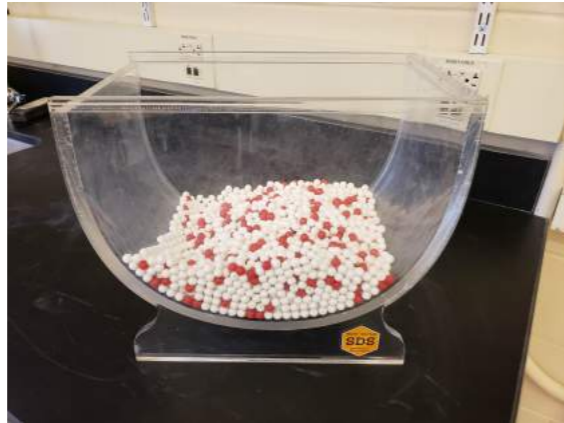
Later: Inference for Regression

1. What are we doing ❓
   - *Descriptive* simple linear regression & regression with single categorical x only.
2. Why are we doing this 🤔
   - Multivariate thinking per GAISE guidelines & modeling
3. Our opinions
   - Separate descriptive vs inference so we can introduce it early, not at end of term 🥵
   - `moderndive::get_regression_table()` function has CI's, no p-value 🌟s
   - Much of world's data is categorical, to skip is to do students a disservice
   - Introduce causal inference
4. Potential pitfalls ⚠️
   - Understanding regression with categorical x

# Chapter 7: Multiple regression
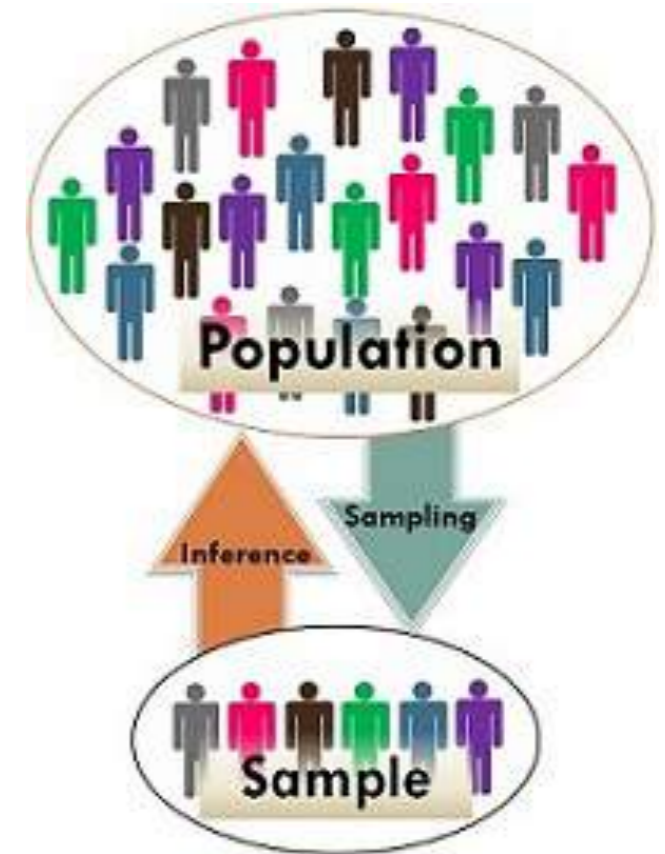
# Goal 1: Sampling for Inference

1. Tactile Sampling → 2. Virtual Sampling → 3. Theoretical

**Population**



```
Console ~/
> library(moderndive)
> bowl
# A tibble: 2,400 x 2
   ball_ID color
     <int> <chr>
1        1 white
2        2 white
3        3 white
4        4 red
5        5 white
6        6 white
7        7 red
8        8 white
9        9 red
10      10 white
# ... with 2,390 more rows
>
```
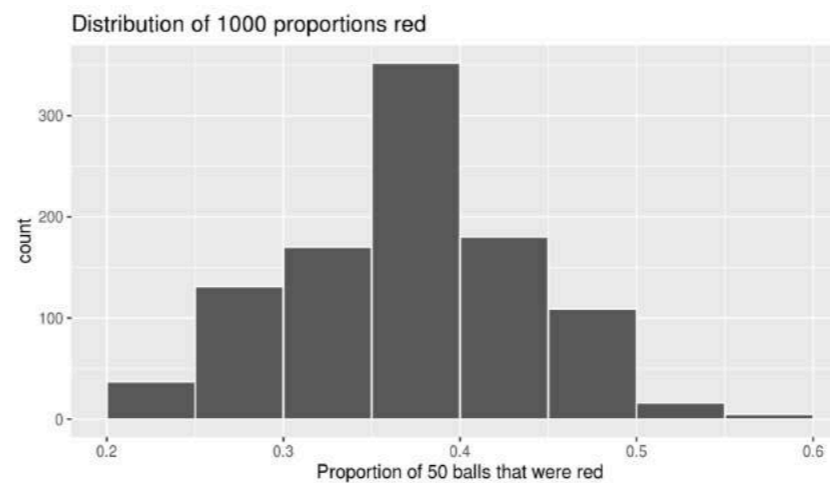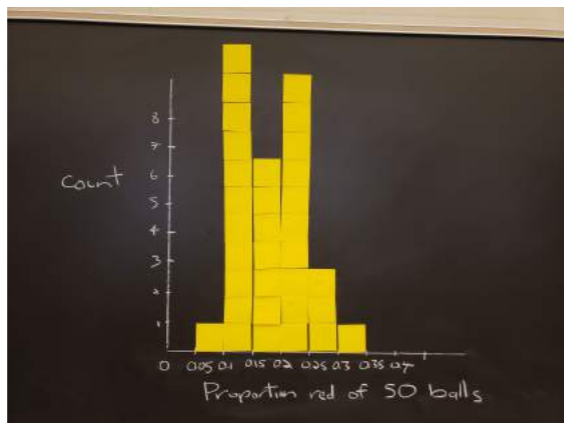
**Sample**



```
Console ~/
> bowl %>%
+   rep_sample_n(size = 50, reps = 1)
# A tibble: 50 x 3
# Groups:   replicate [1]
   replicate ball_ID color
       <int>   <int> <chr>
1          1     226 white
2          1    1304 red
3          1    1230 white
4          1     984 white
5          1      68 white
6          1    1965 white
7          1     431 white
8          1    1184 white
9          1    1610 red
10         1     978 white
# ... with 40 more rows
>
```
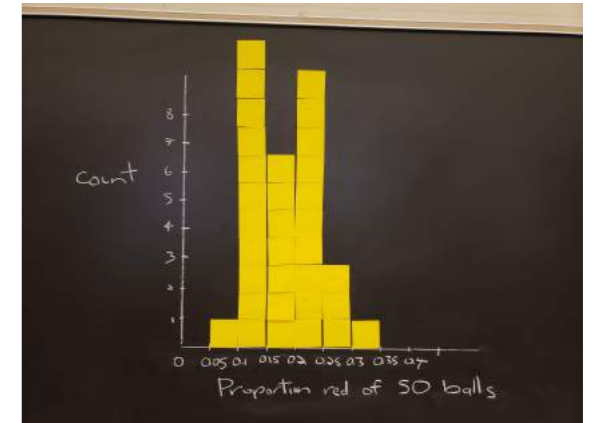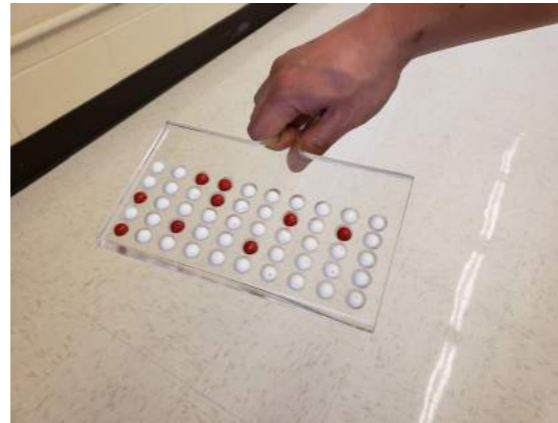
**Sampling Distributions & Standard Errors**



Distribution of 1000 proportions red

$$SE = \sqrt{\frac{p(1-p)}{n}}$$

27

# Statistical Inference



Issues:

1. Easily lost in many layers of abstraction
2. 😱 Notation, terminology, & definitions 😱

# Hypothesis Testing

😱

# Transference & Generalizability

TABLE 8.6: Scenarios of sampling for inference

| Scenario | Population parameter | Notation | Point estimate | Notation. |
|---|---|---|---|---|
| 1 | Population proportion | $p$ | Sample proportion | $\hat{p}$ |
| 2 | Population mean | $\mu$ | Sample mean | $\hat{\mu}$ or $\bar{x}$ |
| 3 | Difference in population proportions | $p_1 - p_2$ | Difference in sample proportions | $\hat{p}_1 - \hat{p}_2$ |
| 4 | Difference in population means | $\mu_1 - \mu_2$ | Difference in sample means | $\bar{x}_1 - \bar{x}_2$ |
| 5 | Population regression slope | $\beta_1$ | Sample regression slope | $\hat{\beta}_1$ or $b_1$ |
| 6 | Population regression intercept | $\beta_0$ | Sample regression intercept | $\hat{\beta}_0$ or $b_0$ |

# Timeline

- **Now:** Follow development version "under construction" at moderndive.netlify.com
- **Mid-June**: Preview of print edition available on moderndive.com
- **Late-July**: Posting labs and example final project samples
- **Fall 2019**: Print edition available!

# Tips!

- Start small
- Teaching/learning code is a psychological battle
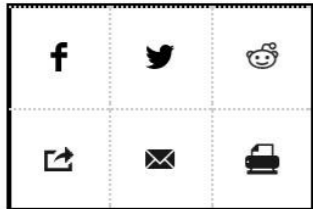- You will run into problems

# The debate continues…

## A Learning Secret: Don't Take Notes with a Laptop

Students who used longhand remembered more and had a deeper understanding of the material

By Cindi May on June 3, 2014     27     Véalo en español



READ THIS NEXT

The Science of Education: Back to School

The old fashioned way works better. *Credit: Credit: Szepy via iStock*

33

# Coding

[Cobb (2015)](#) argued there are two possible computational engines for statistics:

$$t = \frac{(\overline{X}_1 - \overline{X}_2) - (\mu_1 - \mu_2)}{S_{\overline{X}_1 - \overline{X}_2}} = \frac{\overline{X}_1 - \overline{X}_2}{S_{\overline{X}_1 - \overline{X}_2}}$$

$$S_{\overline{X}_1 - \overline{X}_2} = \sqrt{\frac{(N_1 - 1)s_1^2 + (N_2 - 1)s_2^2}{N_1 + N_2 - 2}\left[\frac{1}{N_1} + \frac{1}{N_2}\right]}$$

# Teaching/Learning Code

- Learn how a practitioner would learn: the "Copy/paste/tweak approach"
- Borrow elements of "flipped classroom": how to use time we're all in the same room together?

# Teaching Coding: The Battle is Psychological

- "Don't code from scratch, take the copy/paste/tweak approach!"
- "Computers are stupid!"
- "Learning to code is similar to learning a language"

# New Tools Specific for Data Science

**David Robinson**

*Data Scientist at Stack Overflow, works in R and Python.*

## Teach the tidyverse to beginners

A few years ago, I wrote a post <u>Don't teach built-in plotting to beginners (teach ggplot2)</u>. I argued that ggplot2 was not an advanced approach meant for experts, but rather a suitable introduction to data visualization.
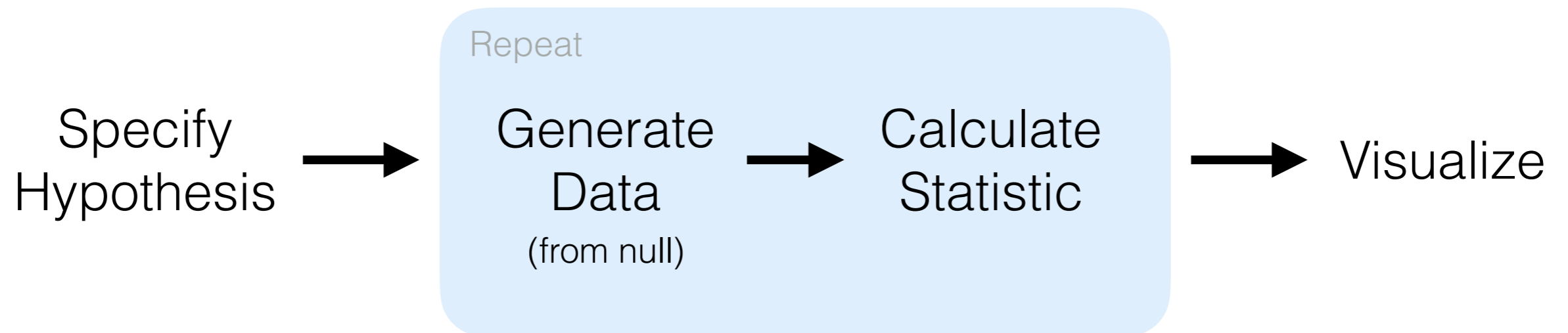
*Many teachers suggest I'm overestimating their students: "No, see, my students are beginners…". If I push the point, they might insist I'm not understanding just how much of a beginner these students are, and emphasize they're looking to keep it simple and teach the basics, and that that students can get to the advanced methods later….*

## DataCamp

# `infer` package for tidy statistical inference

http://infer.netlify.com/



```
hypothesize(null) %>% generate(reps) %>% calculate(stat) %>% visualize()
```

# Hypothesis Testing



**data**

$H_0$

`specify()`

`hypothesize()`

`generate()`

`calculate()`

`visualize()`