# Something old, something new, something borrowed, something **blue**

*Ways to teach data science (and learn it too!)*
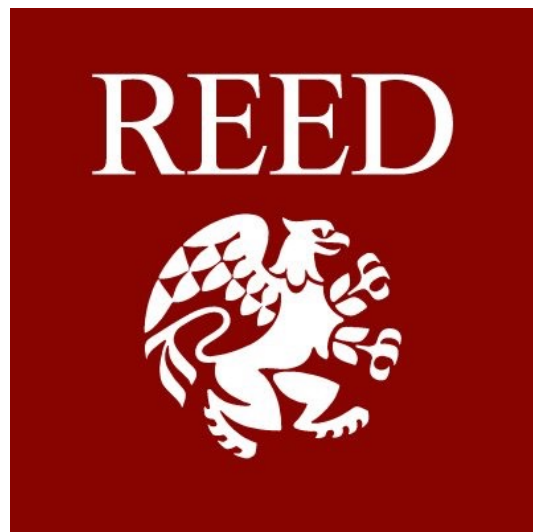


Albert Y. Kim
Amherst College (Smith College July 2018)

Slides available at twitter.com/rudeboybert

# Background

# Focus of Today

- Talk is nominally is about how I teach intro statistics and data science courses
- However can apply to a broader target demographic
- R-centric, but many of these ideas are language agnostic

# Amherst College STAT135

- [Course webpage](#)
- Heterogeneous group: Backgrounds and socio-economics status
- Majors: Math, Stats, Econ, Bio, Neuroscience, Psych, Poli Sci, Environmental Studies
- All had high school algebra, most had no coding experience

# Question

How can we introduce **data and computation** novices to:

1. **Data science**: Data visualization, data wrangling, exploratory data analysis

2. **Data modeling**: Explanation (causal inference) & prediction (machine learning), correlation

3. **Statistical inference**: elementary probability theory, sampling distributions, standard errors, confidence intervals, hypothesis/AB testing & p-values

**An Introduction to Statistical and Data Sciences via R**

- Online textbook available at moderndive.com
- Development version at moderndive.netlify.com
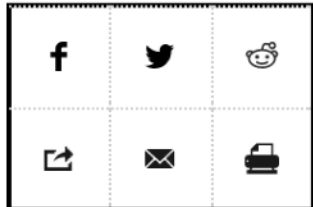- On GitHub at github.com/moderndive/

# Technology in the classroom?

# 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

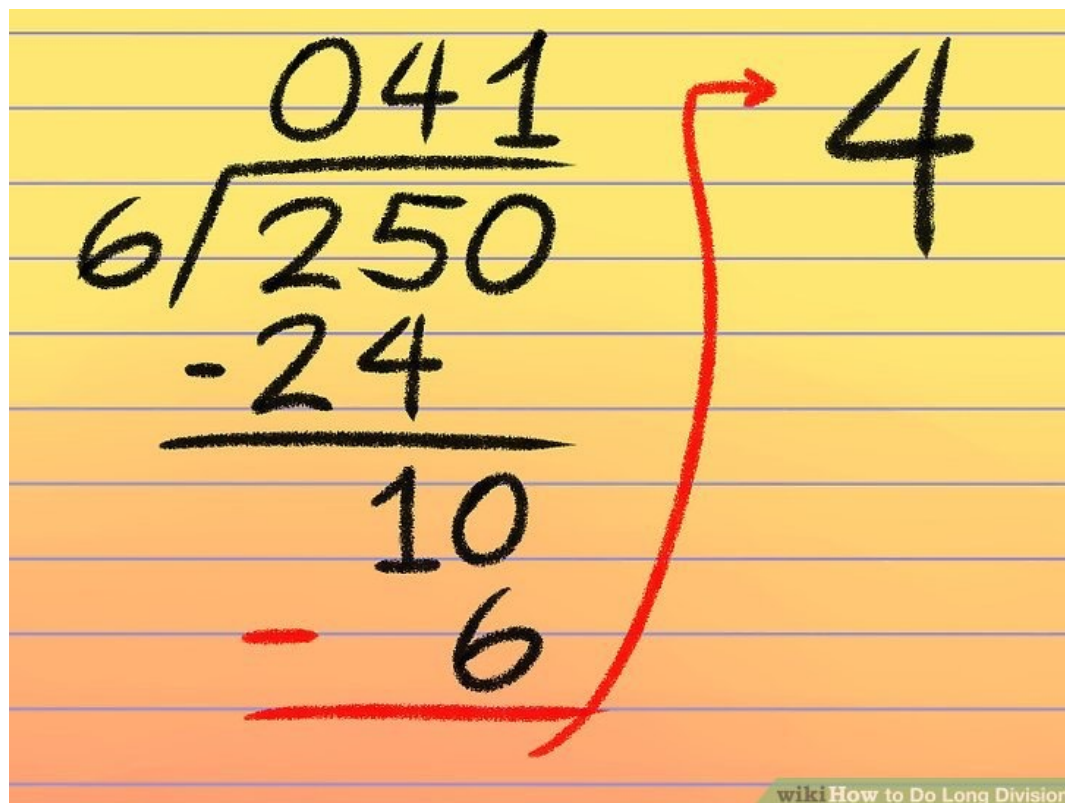By Cindi May on June 3, 2014 · 27 · Véalo en español



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

**READ THIS NEXT**

The Science of Education: Back to School

# Analogy: Learning Long Division

Do this a few times:

Then rely on this:

# `ggplot2` via the Grammar of Graphics

# ggplot2 via the Grammar of Graphics

To create this plot:

① Load ggplot2 package
  library(ggplot2)

② Example of a function call   (problem set 02)
  to create plot in tweet

aesthetic    data variables

$$ggplot(\ data=example,\ aes(x=A,\ y=B,\ color=D\ ))\ +$$

where variables exist

$$geom\_point(\ )$$

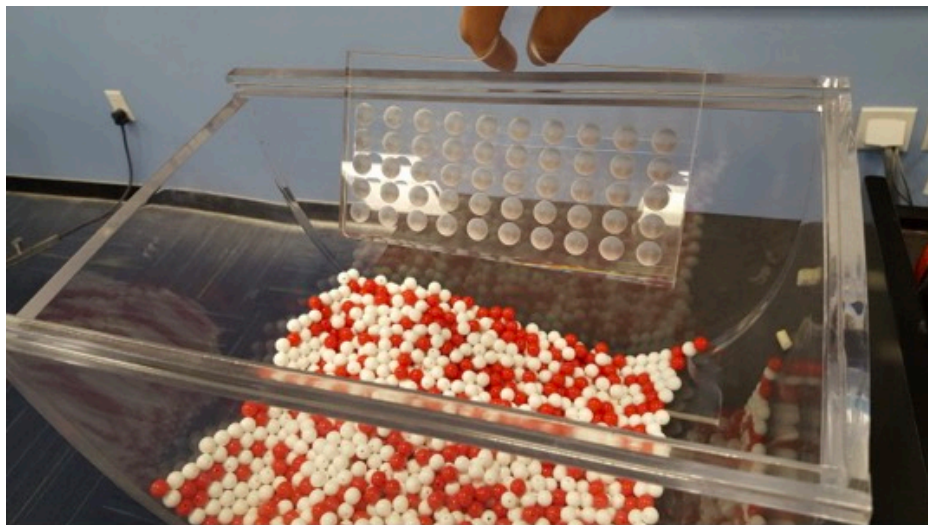geometric object in question.

color

Recall:

A statistical graphic is a mapping of data variables to aesthetic attributes of geometric objects.

Five Named Graphs    5NG
① Scatterplot    geom_point( )
② Linegraphs    geom_line( )
③ Histograms    geom_histogram( )
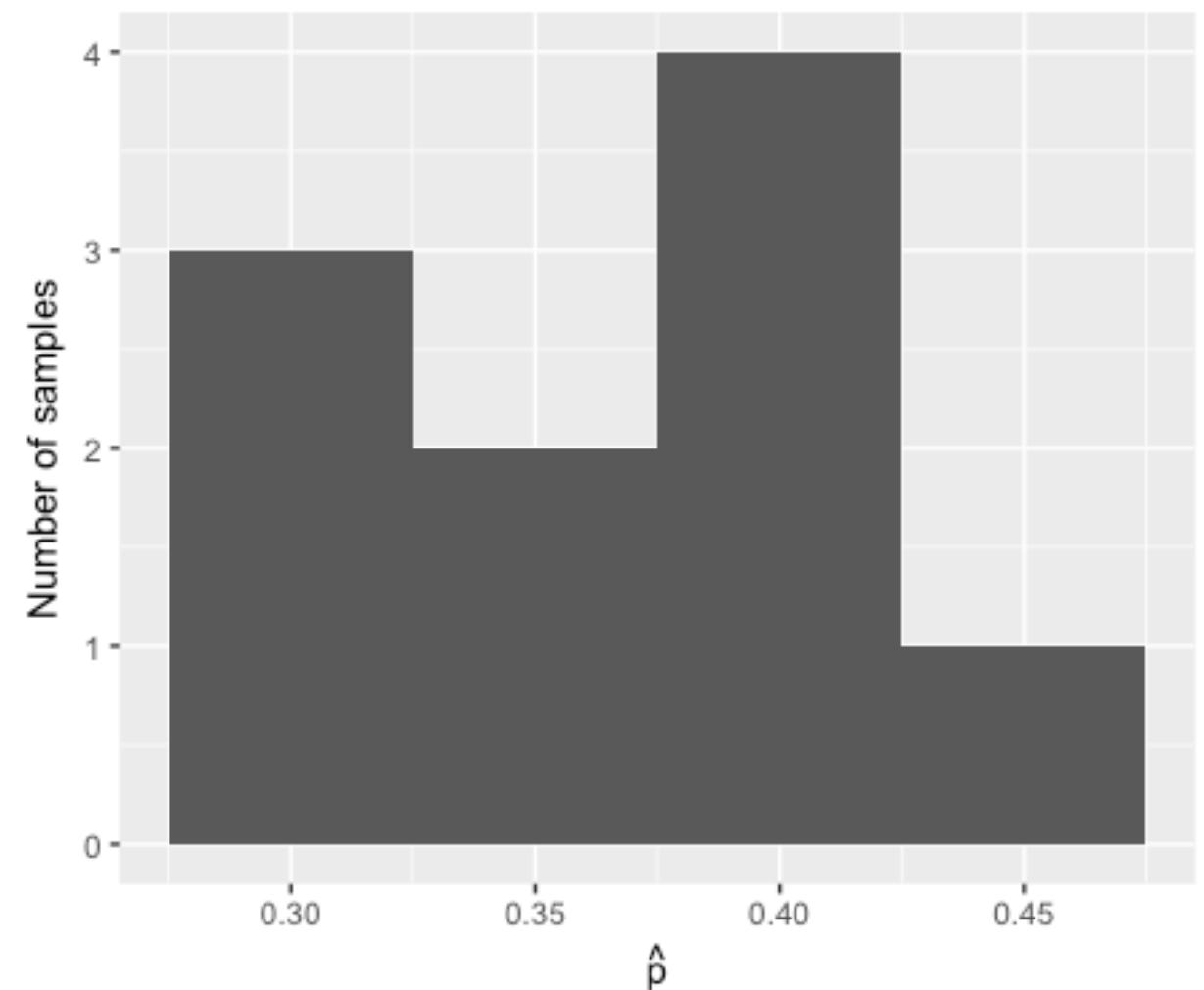④ Boxplots    geom_boxplot( )
⑤ Barplots    geom_bar( )

# Tactile simulation of sampling to teach sampling distributions



| | group | red | n | prop_red |
|---|---|---|---|---|
| 1 | Kathleen and Max | 18 | 50 | 0.36 |
| 2 | Sean, Jack, and CJ | 18 | 50 | 0.36 |
| 3 | X and Judy | 22 | 50 | 0.44 |
| 4 | James and Jacob | 21 | 50 | 0.42 |
| 5 | Hannah and Siya | 16 | 50 | 0.32 |
| 6 | Niko, Sophie, and Caitlin | 14 | 50 | 0.28 |
| 7 | Niko, Sophie, and Caitlin | 19 | 50 | 0.38 |
| 8 | Aleja and Ray | 20 | 50 | 0.40 |
| 9 | Yaw and Drew | 16 | 50 | 0.32 |
| 10 | Yaw and Drew | 21 | 50 | 0.42 |



Sampling distribution of p_hat based on n = 50

# Computer simulation of sampling to teach sampling distributions

| | replicate | red | n | prop_red |
|---|---|---|---|---|
| 1 | 1 | 18 | 50 | 0.36 |
| 2 | 2 | 16 | 50 | 0.32 |
| 3 | 3 | 18 | 50 | 0.36 |
| 4 | 4 | 16 | 50 | 0.32 |
| 5 | 5 | 18 | 50 | 0.36 |
| 6 | 6 | 24 | 50 | 0.48 |
| 7 | 7 | 17 | 50 | 0.34 |
| 8 | 8 | 15 | 50 | 0.30 |
| 9 | 9 | 16 | 50 | 0.32 |
| 10 | 10 | 18 | 50 | 0.36 |

Showing 1 to 10 of 10,000 entries

```
> 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
> bowl %>%
    rep_sample_n(size = 50, reps = 10000)
```



Sampling distribution of p_hat based on n = 50

Something

Old

New

Borrowed

Blue

# 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 Don't teach built-in plotting to beginners (teach ggplot2). 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: Immediate Feedback

- Students can practice failing, but with support.
- Difference with Coursera & Udacity?
- DataCamp will pick off low hanging fruit. Ex:
  1. Matching parentheses
  2. Variable name misspellings
  3. Linearity of programs
- Examples of "Curse of knowledge"

# **Without** DataCamp: # of Questions on Coding

# **With** DataCamp: # of Questions on Coding

# Leverage open source

Open data, such as data in R packages like
nycflights13, gapminder, `fivethirtyeight`

Bechdel test?   Original 538 article

# Leverage open source

# New textbook authoring paradigm

# New textbook authoring paradigm
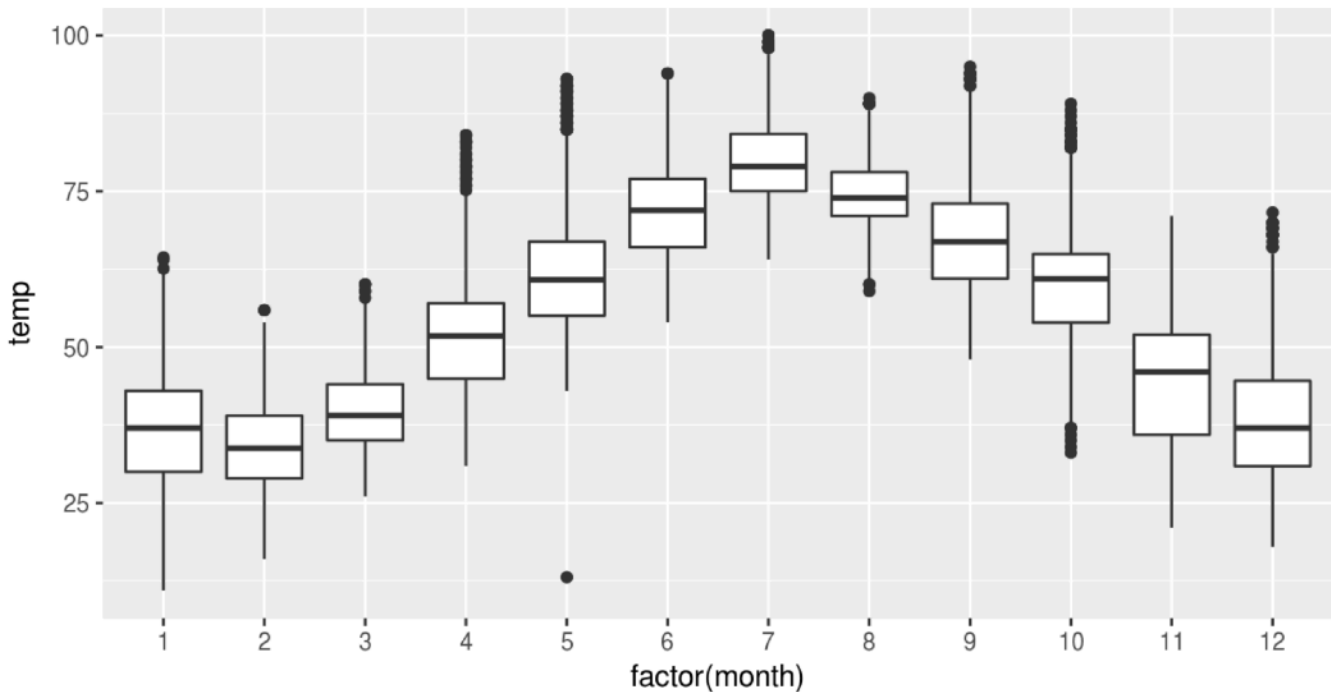
# New textbook authoring paradigm



Figure 3.13: Month by temp boxplot

We have introduced a new function called `factor()` here. One of the things this function does is to convert a discrete value like `month` (1, 2, ..., 12) into a categorical variable. The "box" part of this plot represents the 25[th] percentile, the median (50[th] percentile), and the 75[th] percentile. The dots correspond to *outliers*. (The specific formulation for these outliers is discussed in Appendix A.) The lines show how the data varies that is not in the center 50% defined by the first and third quantiles. Longer lines correspond to more variability and shorter lines correspond to less variability.

# New textbook authoring paradigm



**"Versions, not editions"**

On GitHub at github.com/moderndive/

**An Introduction to Statistical and Data Sciences via R**

- Available at moderndive.com
- Development version at moderndive.netlify.com
- On GitHub at github.com/moderndive/

**v0.3.0 to be released next week! What's new?**

1. Introduction

2. Getting Started
with Data in R

3. Data
Visualization

4. Tidy
Data

5. Data
Wrangling

**Data Science with** `tidyverse`

ModernDive

**Available at** `moderndive.com`

**Diagram inspired by hadley/r4ds**

12. Thinking with Data

6. Basic
Regression

11. Inference
for Regression

7. Multiple
Regression

**Data Modeling with** `moderndive`

8. Sampling

9. Confidence
Intervals

10. Hypothesis
Testing

**Statistical Inference with** `infer`

*"If You're Not Embarrassed By The First Version Of Your Product, You've Launched Too Late"*

[Reid Hoffman, founder of LinkedIn](#)

# Crowdsourcing Typos

# `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()**
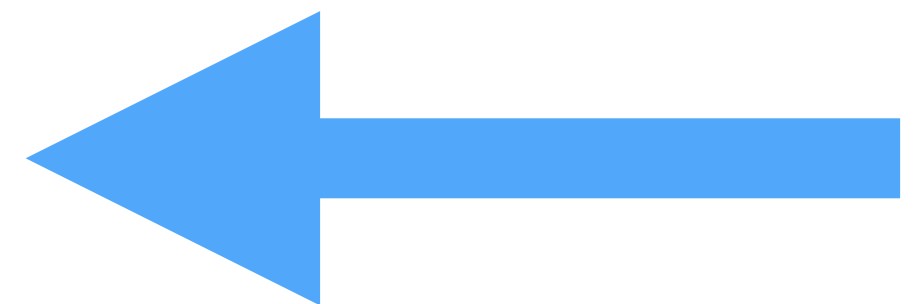
# "Thinking with Data"

Example student work

- Analysis of crime in [Chicago](#)
- How many [f**ks](#) does Tarantino Give?
- Final projects: [Code and data](#)

RStudio + R + ModernDive

Albert Y. Kim
Amherst College
Twitter: @rudeboybert
GitHub: rudeboybert

Chester Ismay
DataCamp
Twitter: @old_man_chester
GitHub: ismayc

**An Introduction to Statistical and Data Sciences via R**

- Available at moderndive.com
- Development version at moderndive.netlify.com
- On GitHub at github.com/moderndive/

**v0.3.0 to be released next week! What's new?**

Slides available at twitter.com/rudeboybert