

expanding   
exposure through  
early introduction  
in the undergraduate  
curriculum

mine çetinkaya-rundel  
duke university





intro stats  
(general ed)

extra-  
curriculars

stats major  
+ other  
quantitative  
fields

intro stats  
(general ed)

why R?

why not R?

unlike other software designed  
specifically for courses at this level

free & open  
source

powerful &  
flexible

relevant  
beyond intro  
stat

challenge of  
teaching  
programming in  
addition to stats  
concepts

command line  
more  
intimidating  
than GUI

challenge of  
teaching  
programming in  
addition to stats  
concepts

don't do any  
hands on data  
analysis

disservice to  
everyone  
involved



challenge of  
teaching  
programming in  
addition to stats  
concepts

use a drag-and-  
drop type tool

don't do any  
hands on data  
analysis

there's still a  
learning  
curve

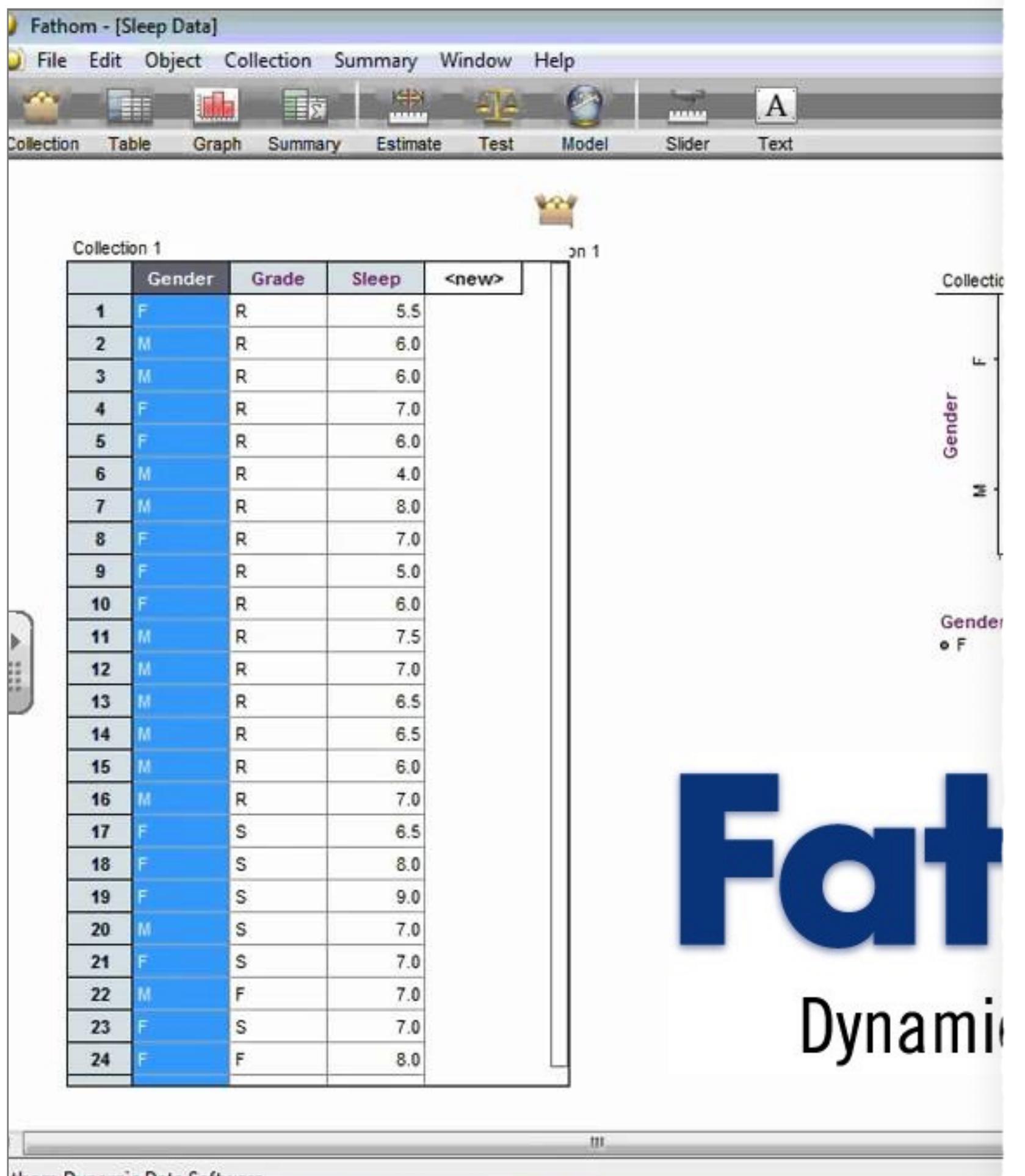
disservice to  
everyone  
involved



### **III. Adding Proportions to Summary Table**

For categorical variables, you should see the counts of each possible outcome of that variable in the **Summary Table**. To see the breakdown of proportions or percentages, follow these steps:

- Click on the **Summary Table** to highlight it, click on the “**Summary**” drop-down menu and select “**Add Formula**”. In general, whenever you click and select a *Fathom* object (such as a **Table**, **Graph**, or **Summary**) the menu at the top of the screen will change to give you options for working on that object.
- In the formula editor that pops up, type “*rowproportion*” (without the quotes) to see the row proportions or “*columnproportion*” to see the column proportions. Be sure to spell the names of the formulas correctly or else *Fathom* will give you an error. (If you spell the names correctly, they should change to a purplish color in your editor.)
- You will see that each cell in the **Summary Table** now includes numbers for multiple statistics. To see which numbers correspond with which statistics, simply look at the bottom of your summary table to see the order of the statistics or formulas within each cell.
- To delete (or change) a particular statistic from the table, you can double click on its name at the bottom of the **Summary Table**. In the formula editor, press delete (or make your changes) and then click “**OK**”.



challenge of  
teaching  
programming in  
addition to stats  
concepts

use a drag-and-  
drop type tool

don't do any  
hands on data  
analysis

there's still a  
learning  
curve

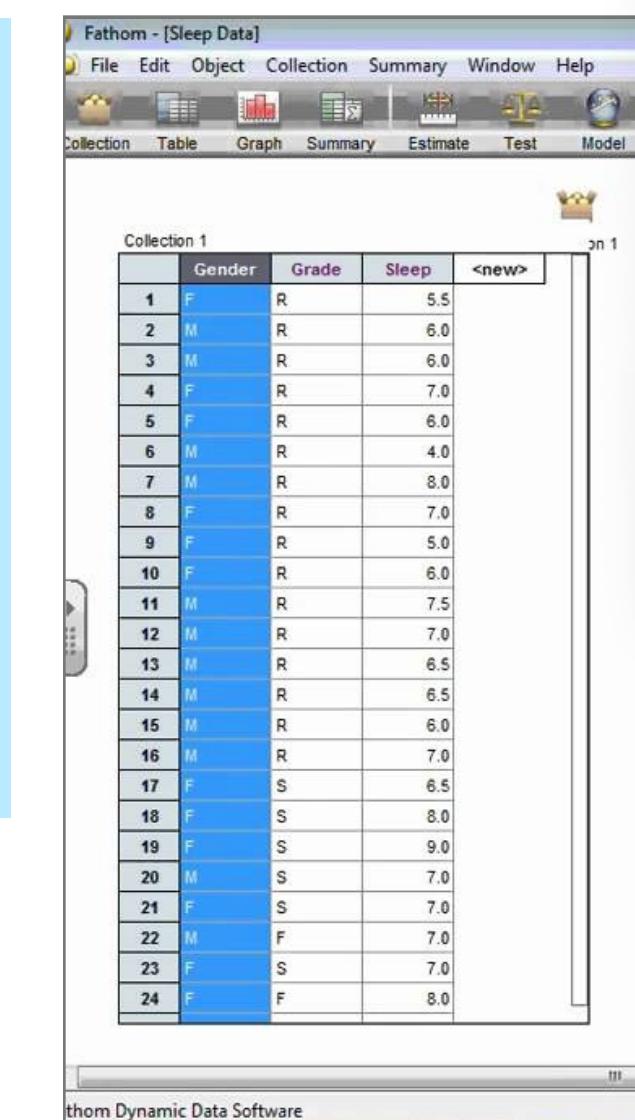
disservice to  
everyone  
involved



### III. Adding Proportions to Summary Table

For categorical variables, you should see the counts of each possible outcome of that variable in the **Summary Table**. To see the breakdown of proportions or percentages, follow these steps:

- Click on the **Summary Table** to highlight it, click on the “Summary” drop-down menu and select “Add Formula”. In general, whenever you click and select a *Fathom* object (such as a **Table**, **Graph**, or **Summary**) the menu at the top of the screen will change to give you options for working on that object.
- In the formula editor that pops up, type “*rowproportion*” (without the quotes) to see the row proportions or “*columnproportion*” to see the column proportions. Be sure to spell the names of the formulas correctly or else *Fathom* will give you an error. (If you spell the names correctly, they should change to a purplish color in your editor.)
- You will see that each cell in the **Summary Table** now includes numbers for multiple statistics. To see which numbers correspond with which statistics, simply look at the bottom of your summary table to see the order of the statistics or formulas within each cell.
- To delete (or change) a particular statistic from the table, you can double click on its name at the bottom of the **Summary Table**. In the formula editor, press delete (or make your changes) and then click “OK”.



**Fathom®**  
Dynamic Data Software

command line  
more  
intimidating  
than GUI

R version 3.2.1 (2015-06-18) -- "World-Famous Astronaut"  
Copyright (C) 2015 The R Foundation for Statistical  
Computing  
Platform: x86\_64-apple-darwin13.4.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.  
You are welcome to redistribute it under certain conditions.  
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

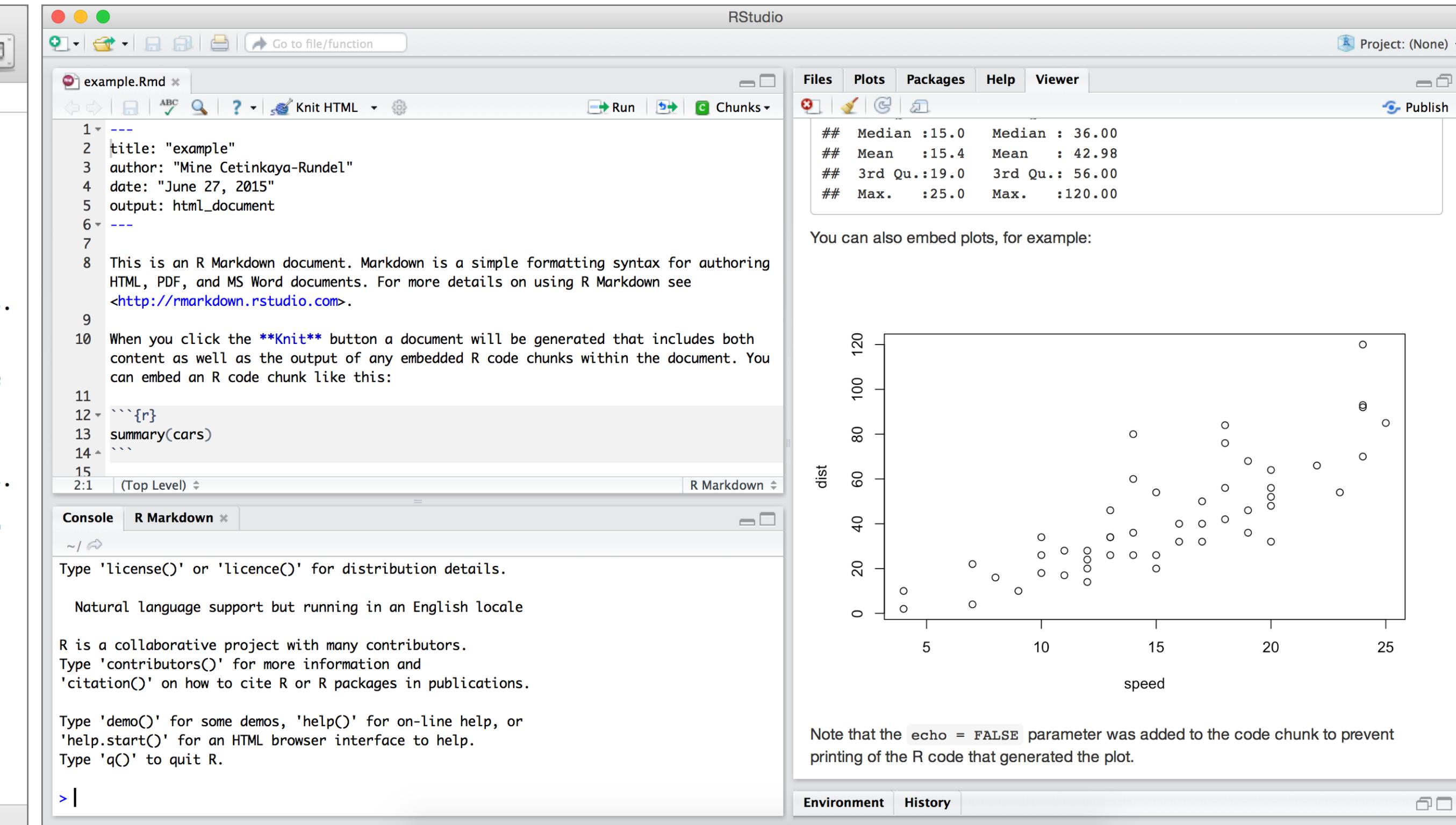
R is a collaborative project with many contributors.  
Type 'contributors()' for more information and  
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.

[R.app GUI 1.66 (6956) x86\_64-apple-darwin13.4.0]

[History restored from /Users/mine/.Rhistory]

> |



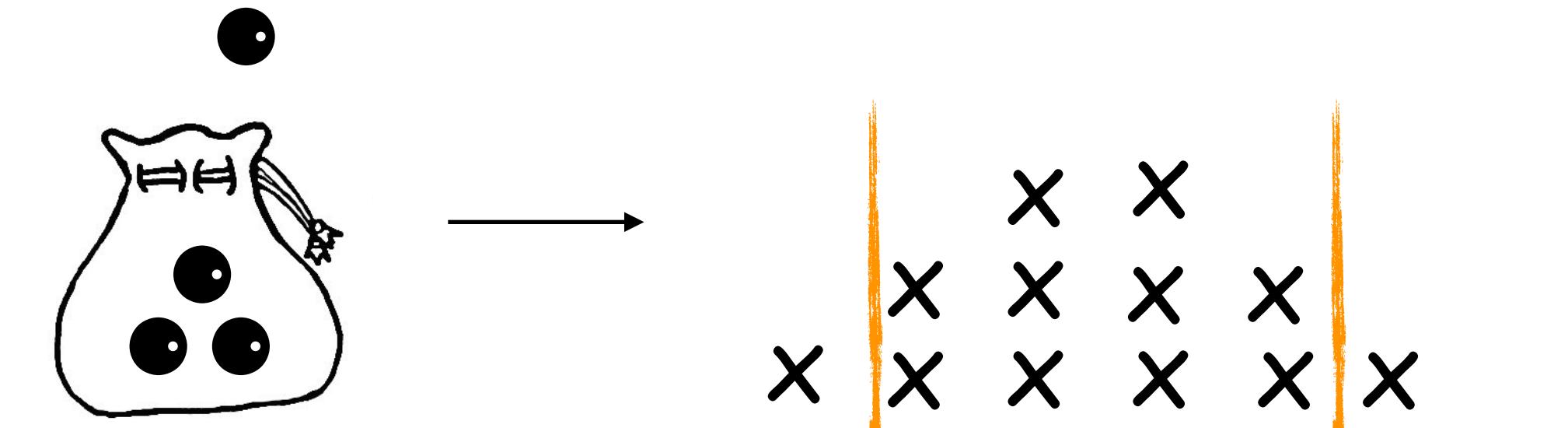
# how R?

**balance:** teach code as a way of introducing/reinforcing concepts

bootstrapping

physical representation

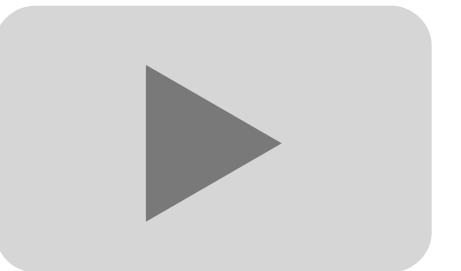
computation



```
for(i in 1:15000){  
  samp <- sample(mydata, size = n)  
  boot_dist[i] <- mean(mydata)  
}  
hist(boot_dist)
```

# how R?

**balance:**  
minimize  
code for  
implementation  
outside the  
scope of course



# how R?

computation to  
help solve real  
world problems

visualization  
of multiple  
variables at  
once

cleaning of  
real not-very-  
small data

modeling

...

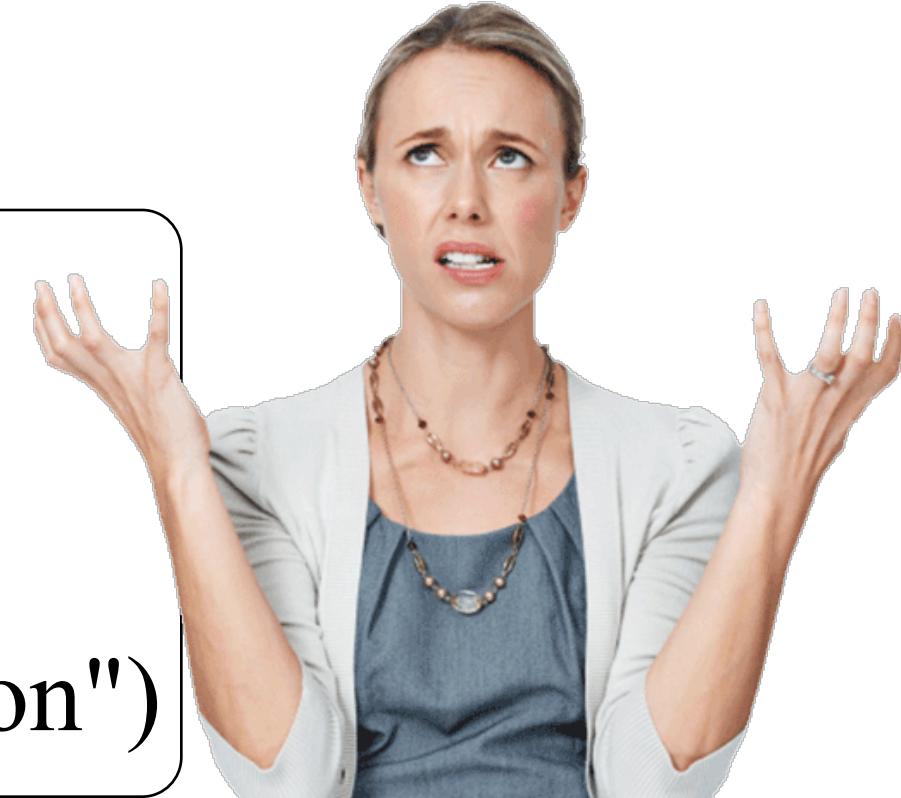
# how R?

consistent  
visual syntax  
highlighting

```
n <- 1000
p <- seq(0, 1, 0.01)
me <- 2 * sqrt(p * (1 - p)/n)
plot(me ~ p, ylab = "Margin of Error", xlab = "Population Proportion")
```



```
n <- 1000
p <- seq(0, 1, 0.01)
me <- 2 * sqrt(p * (1 - p)/n)
plot(me ~ p, ylab = "Margin of Error", xlab = "Population Proportion")
```



# how R?

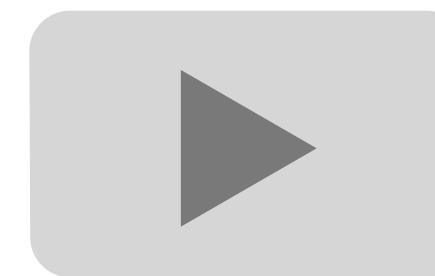
as little  
overhead as  
possible

avoid local  
installation

RStudio  
server  
(via university  
login)



preinstalled &  
preloaded  
packages



**reproducible**  
via literate  
programming

how R?

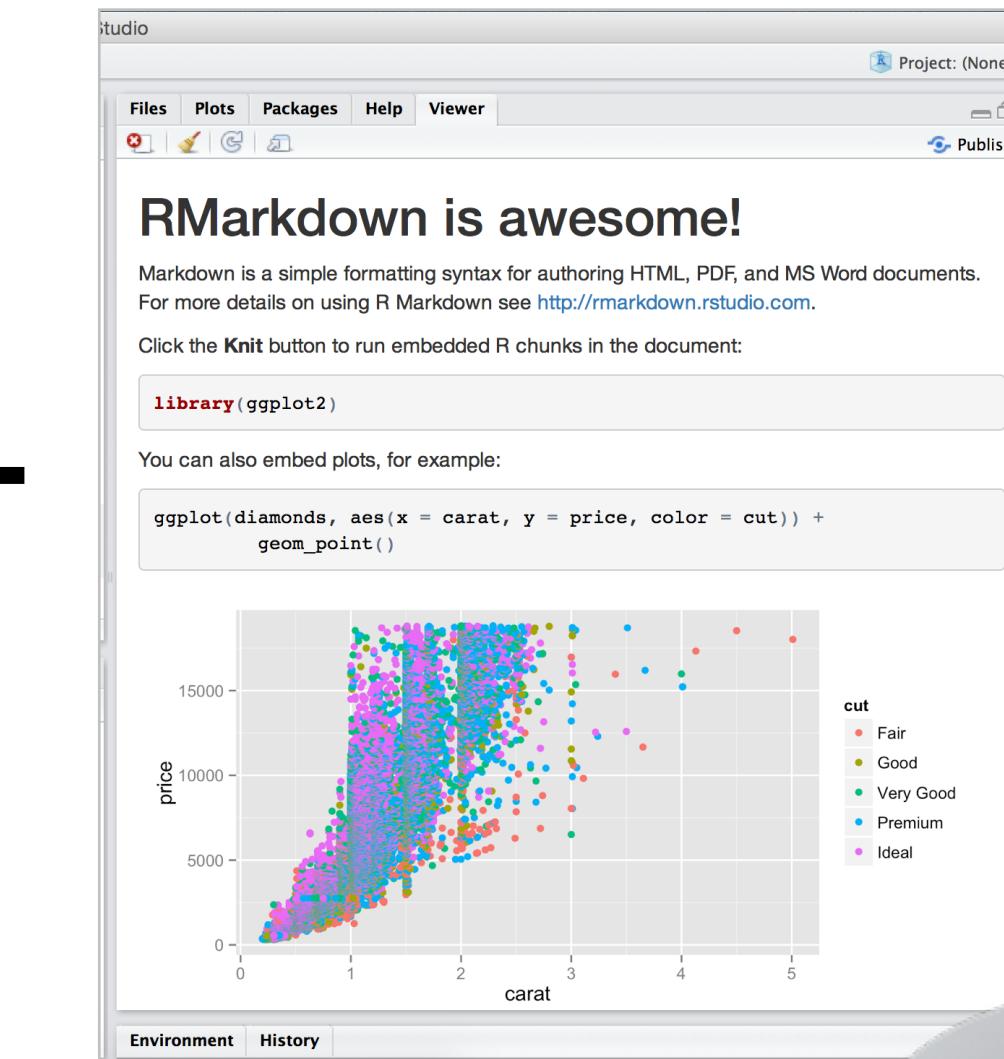
**goal:** train new  
researchers  
whose only  
workflow is a  
reproducible  
one

don't touch  
the raw data

keep track of  
all analysis  
steps

avoid copy-  
paste

**toolkit:**



= Literate programming in

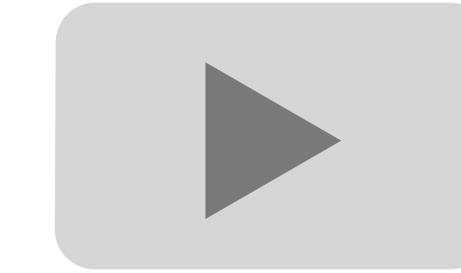


# how R?

with lots of  
scaffolding

start with  
templates  
including code  
and answers

slowly remove  
handholding



R Markdown  
learning  
outcomes  
(beyond  
reproducibility)

learn R

built-in and  
consistent  
syntax  
highlighting

code and  
output always  
together

avoid the  
messy /  
frustrating  
console

feedback +  
grading

ambiguity  
removed

collaboration

Room for  
improvement



intro stats  
(general ed)

extra-  
curriculars

stats major  
+ other  
quantitative  
fields

# ASA DataFest™

weekend-long  
data analysis  
competition

not an R  
event, but R is  
most popular

pre-event  
workshops

## WORKSHOPS

*These workshops are recommended for DataFest participants, but they're open to everyone.  
All workshops will be held in the The Edge Workshop Room.*

» [Intro to R](#)

[slides](#) | [source](#)

Monday, March 16, 6-8pm

**Instructors:** Gary Larson and Monika Hu, Duke StatSci

**Description:** Introduction to R as a statistical programming language. This session will introduce the basics of R syntax, getting data into R, various data types and classes, etc. The session assumes no or little background in R.

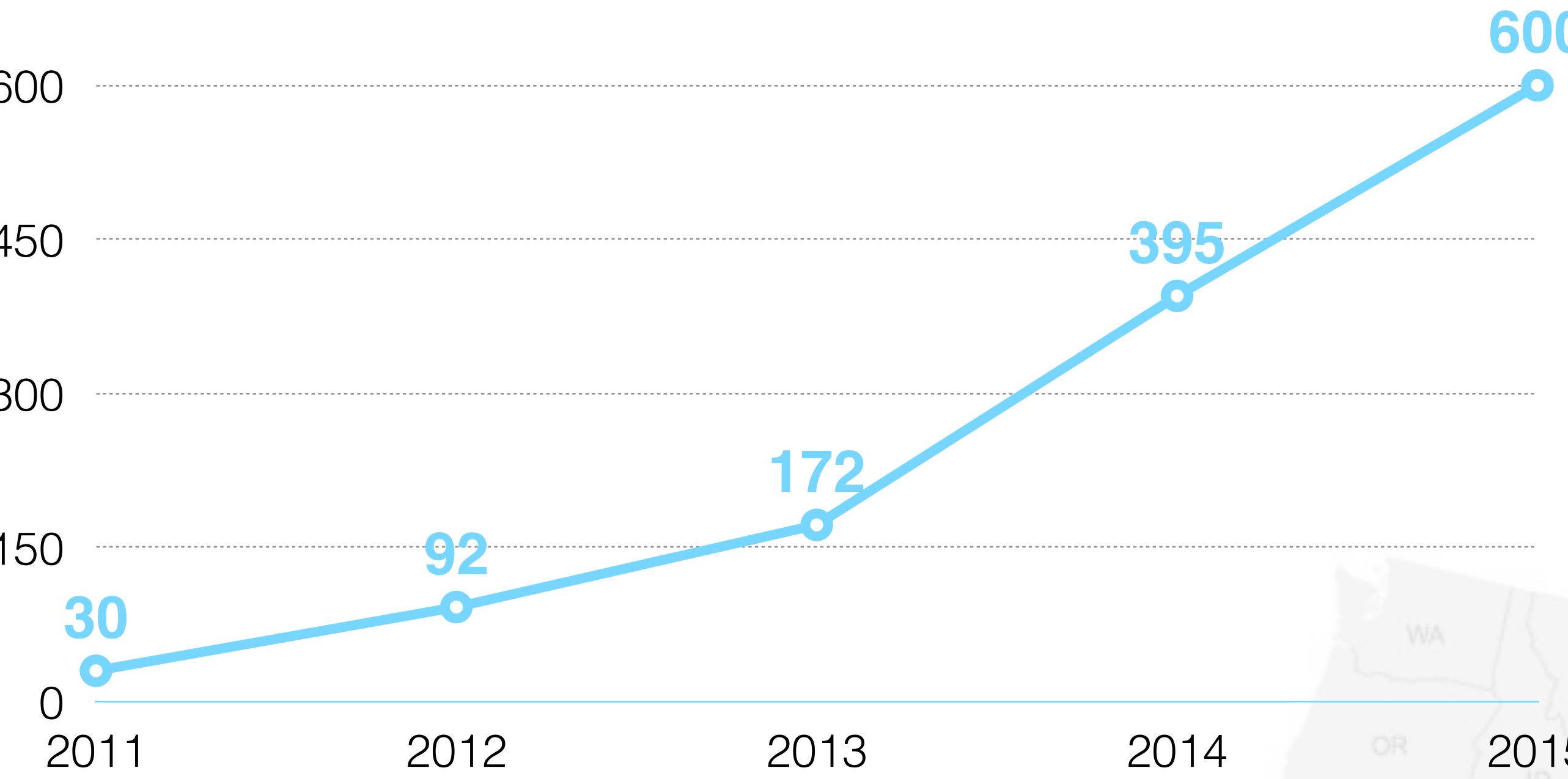
» [Data munging with R and dplyr](#)

[slides](#) | [solutions](#) | [source](#)

Wednesday, March 18, 6-8pm

**Instructor:** Prof. Colin Rundel, Duke StatSci

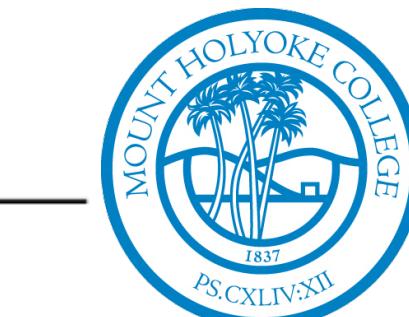
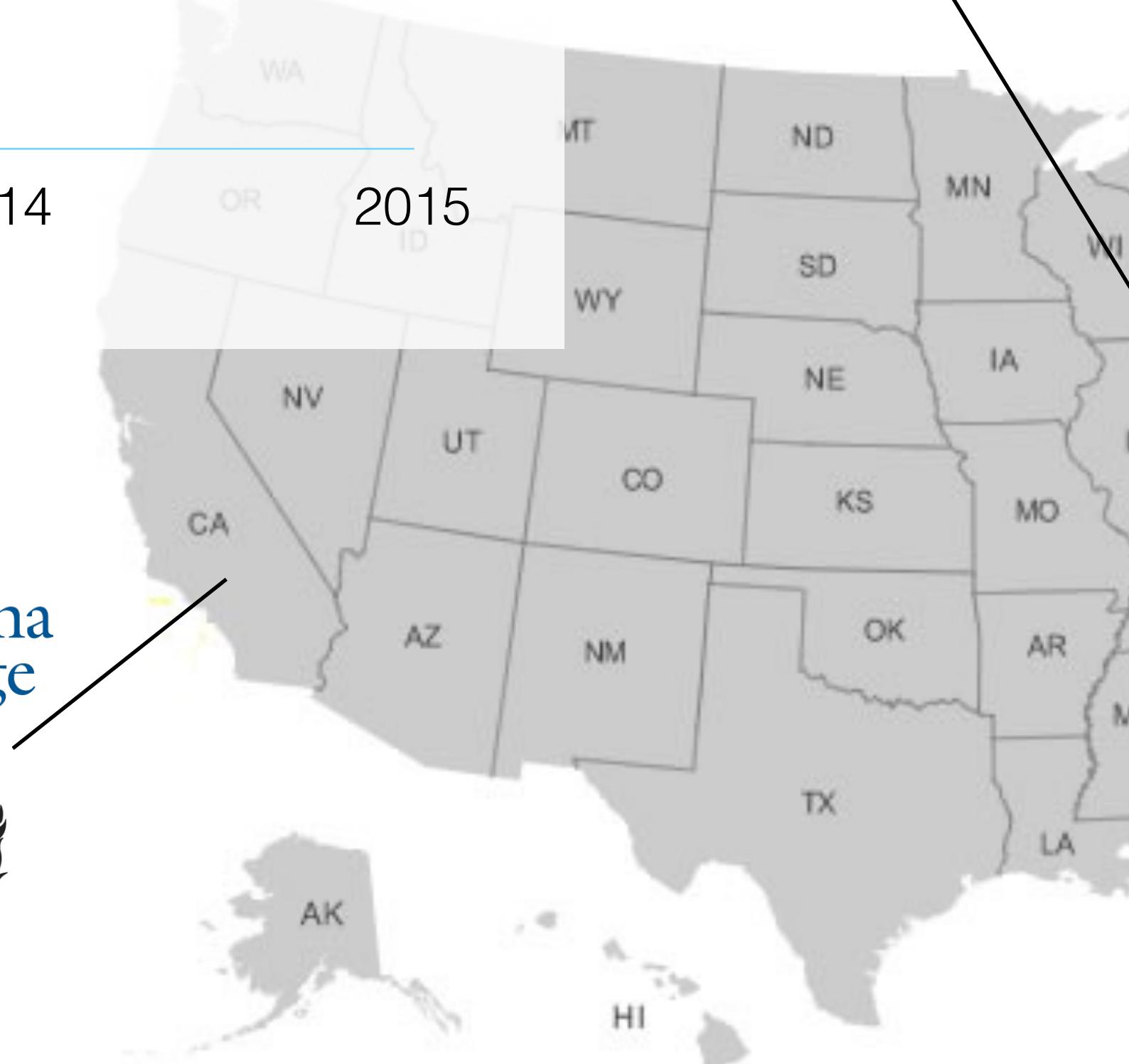
**Description:** This session will demonstrate tools for data manipulation and cleaning of data in R. Majority of the session will use the dplyr and tidyverse packages. Some background in R is recommended.



UNIVERSITY OF CALIFORNIA  
**UCRIVERSIDE**



Pomona  
College



**UMASS  
AMHERST**



THE GEORGE  
WASHINGTON UNIVERSITY  
WASHINGTON, DC



The seal of North Carolina Central University is circular. It features a central figure of a man in academic regalia (gown and cap) standing on a small hill. The background behind him shows a landscape with trees and a path. The text "NORTH CAROLINA CENTRAL UNIVERSITY" is written in a circular pattern around the top edge of the seal. At the bottom, the text "Durham, North Carolina" is written, flanked by two ribbon-like banners that extend upwards. The left banner contains the word "TRUTH" and the right banner contains the word "SERVICE". The year "1910" is positioned at the bottom center between the two banners.



EMORY  
UNIVERSITY



North Carolina  
Agricultural and Technical  
State University

# ASA DataFest™

inform the  
curriculum

learning  
computation  
tends to be  
ad hoc

demand for  
formal stat  
computing  
course

focusing  
especially on  
“the details”



intro stats  
(general ed)

extra-  
curriculars

stats major  
+ other  
quantitative  
fields

stats major  
+ other  
quantitative  
fields

stats elective  
for majors

data science  
course for 1st  
year  
undergrads

stat  
computing  
course for  
students with  
regression

same  
principles as  
intro stat

focus:  
working with  
large &  
complex data

heavier  
emphasis on  
computation

focus: data  
wrangling +  
visualization

# reach

# students, beyond duke

# teachers

# beyond the discipline

# workshops

# open source resources



[openintro.org](http://openintro.org)



# Feedback: Resources for RStudio server setup



thank you!  
comments / questions?



[mine@stat.duke.edu](mailto:mine@stat.duke.edu)



@minebocek



mine-cetinkaya-rundel