

# Introduction to Machine Learning – Evaluating Models

Dr. Ab Mosca (they/them)

Slides based off slides courtesy of Jordan Crouser (<https://jcrouser.github.io/>)

# Plan for Today

- Evaluating:
  - Regression
  - Classification
  - Bias-variance trade off
- R demo time permitting

# What we'll cover in this class

- Ch. 2: Statistical Learning Overview (today)
- Ch. 3: Linear Regression
- Ch. 4: Classification
- Ch. 5: Resampling Methods
- Ch. 6: Linear Model Selection
- Ch. 7: Beyond Linearity
- Ch. 8: Tree-Based Methods
- Ch. 9: Support Vector Machines
- Ch. 10: Unsupervised Learning



Preparing for  
labs in R

You can install R Studio on your own machine:  
[rstudio.com](https://rstudio.com)

# Preparing for labs in python



- You can download the Anaconda distribution from [continuum.io](https://continuum.io) (or a different source)
- You'll need to know how to **install packages**

# Final Project

- Topic: ANYTHING YOU WANT
- Goals:
  - Learn how to break big, unwieldy questions down into clear, manageable problems
  - Figure out if/how the techniques we cover in class apply to your specific problems
  - **Use ML to address them**
- Several (graded) milestones along the way
- Demos and discussion on the final day of class
- More on this later...

One model to  
rule them  
all...?

**Question:** why not just teach you the **best** method first?



Answer: there  
isn't one

- No single method dominates
- One method may prove useful in answering some questions on a given data set
- On a related (not identical) dataset or question, another might prevail





## Measuring “quality of fit”

- *Question we often ask:* how **good** is my model?
- *What we usually mean:* how well do my model's predictions **actually match** the observations?

How do we choose the **right approach**?

# Mean squared error

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

True response  
for the  $i^{th}$  observation

We take the average  
over all observations

Prediction our model gives  
for the  $i^{th}$  observation

The diagram illustrates the components of the Mean Squared Error (MSE) formula. The formula is  $MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$ . Three arrows point to specific parts of the formula: one from the text 'True response for the  $i^{th}$  observation' to  $y_i$ , one from 'We take the average over all observations' to the denominator  $n$ , and one from 'Prediction our model gives for the  $i^{th}$  observation' to  $\hat{y}_i$ .

# Mean squared error

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

True response for the  $i^{th}$  observation

Prediction our model gives for the  $i^{th}$  observation

We take the average over all observations

Prediction our model gives for the  $i^{th}$  observation

| Student | Grade | $f(X)$ |
|---------|-------|--------|
| Ab      | 83    | 78     |
| Kaden   | 84    | 85     |
| Kylee   | 95    | 65     |

$n = 3$

$$\begin{aligned} MSE &= \frac{1}{3} \sum_{i=1}^3 (y_i - \hat{y}_i)^2 \\ &= \frac{1}{3} ((y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + (y_3 - \hat{y}_3)^2) \\ &= \frac{1}{3} ((83 - 78)^2 + (83 - 85)^2 + (95 - 65)^2) \\ &= 309.67 \end{aligned}$$

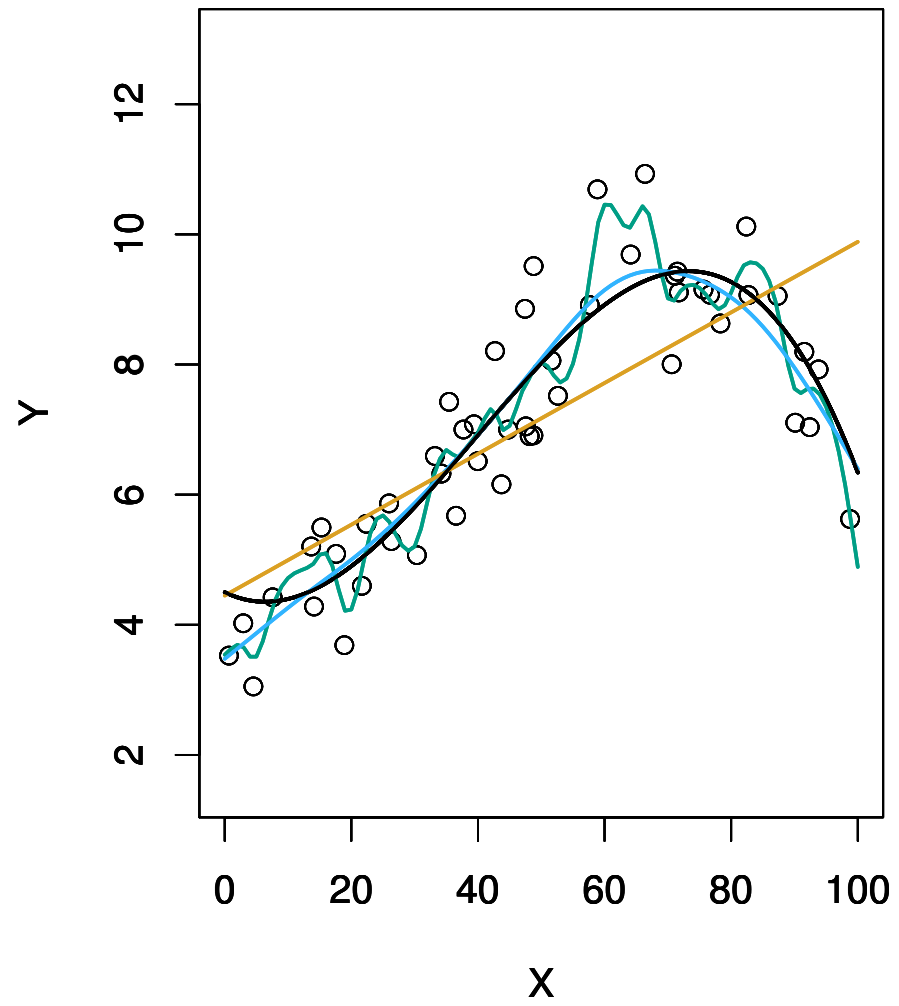
## “Training” MSE

- This version of MSE is computed using the **training data** that was used to fit the model
- **Reality check:** is this what we care about?

# Test MSE

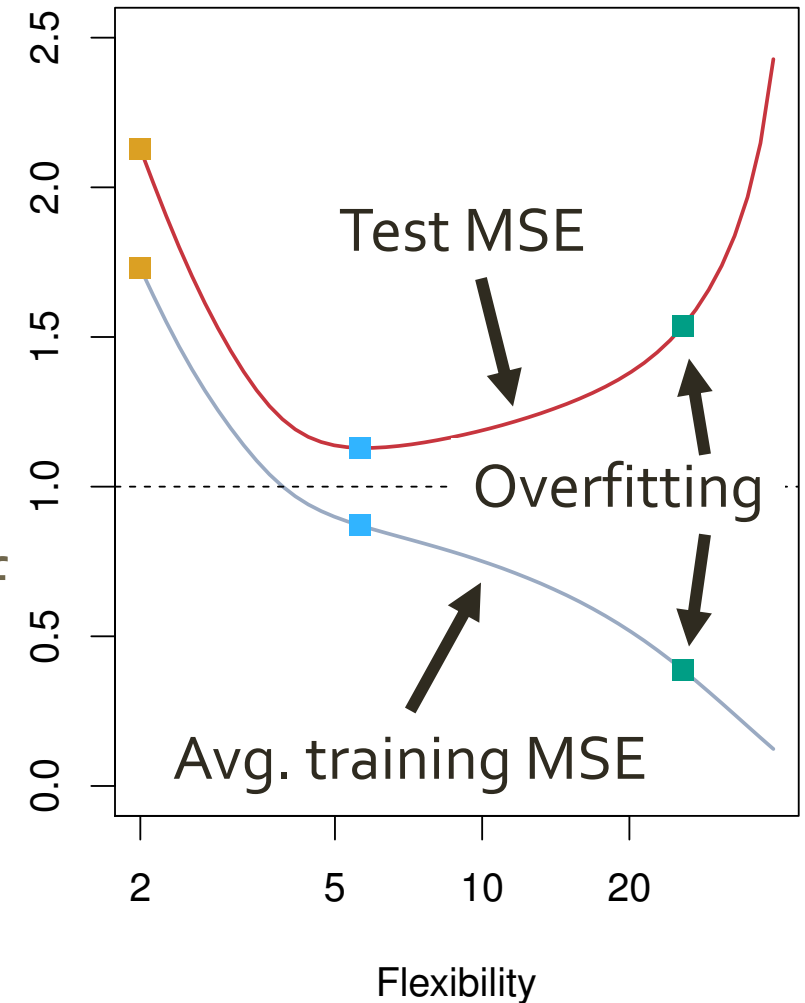
- **Better plan:** see how well the model does on observations we *didn't* train on
- Given some never-before-seen examples, we can just calculate the MSE on those using the same method
- What if we don't have any new observations to test?
  - Can we just use the training MSE?
  - Why or why not?

# Example



# Training vs. test MSE

- As flexibility  $\uparrow$ :
  - monotone  $\downarrow$  in training MSE
  - U-shape in the test MSE
- **Fun fact:** occurs regardless of data or statistical method
- This is called **overfitting**



Training vs.  
test MSE

**Question:** why does this happen?

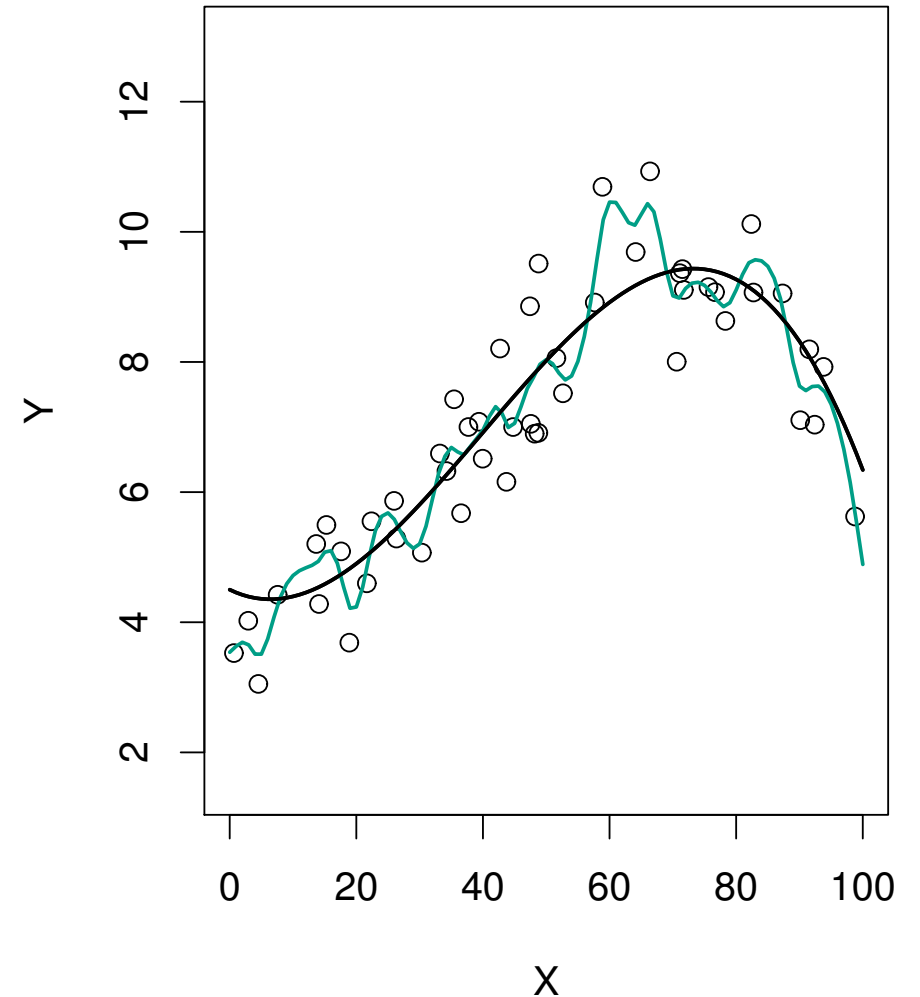
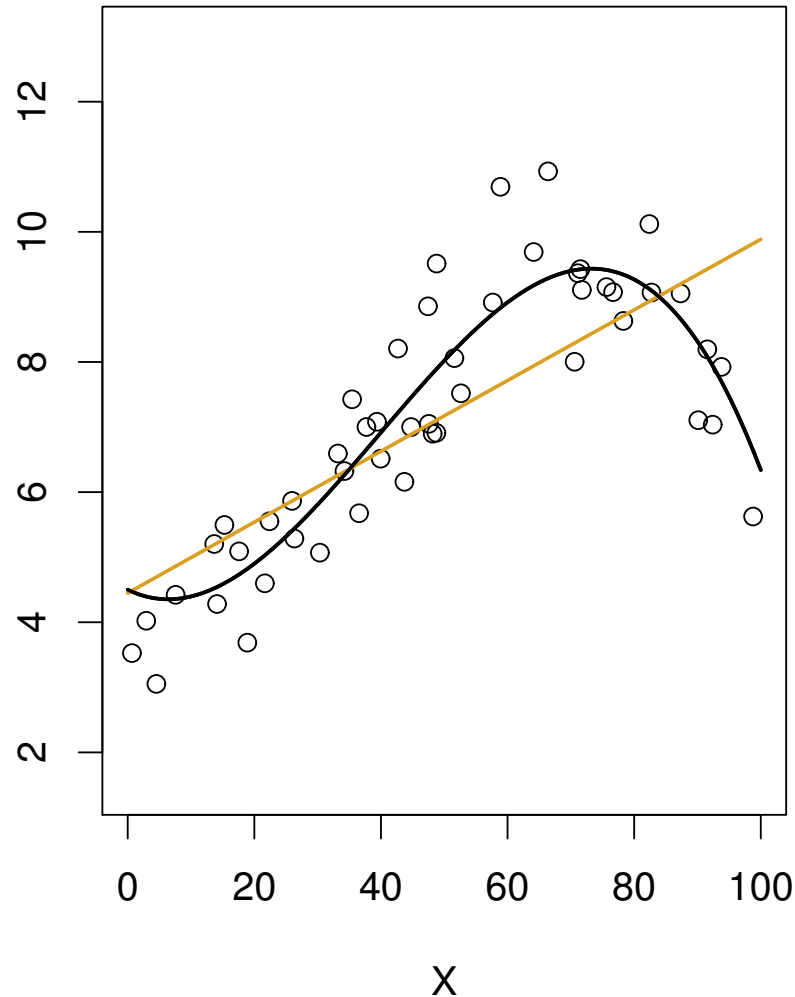


## Trade-off between bias and variance

- The U-shaped curve in the Test MSE is the result of two competing properties: *bias* and *variance*
- **Variance:** the amount the model would change if we had different training data
- **Bias:** the error introduced by approximating a complex phenomenon using a simple model

# Relationship between bias and variance

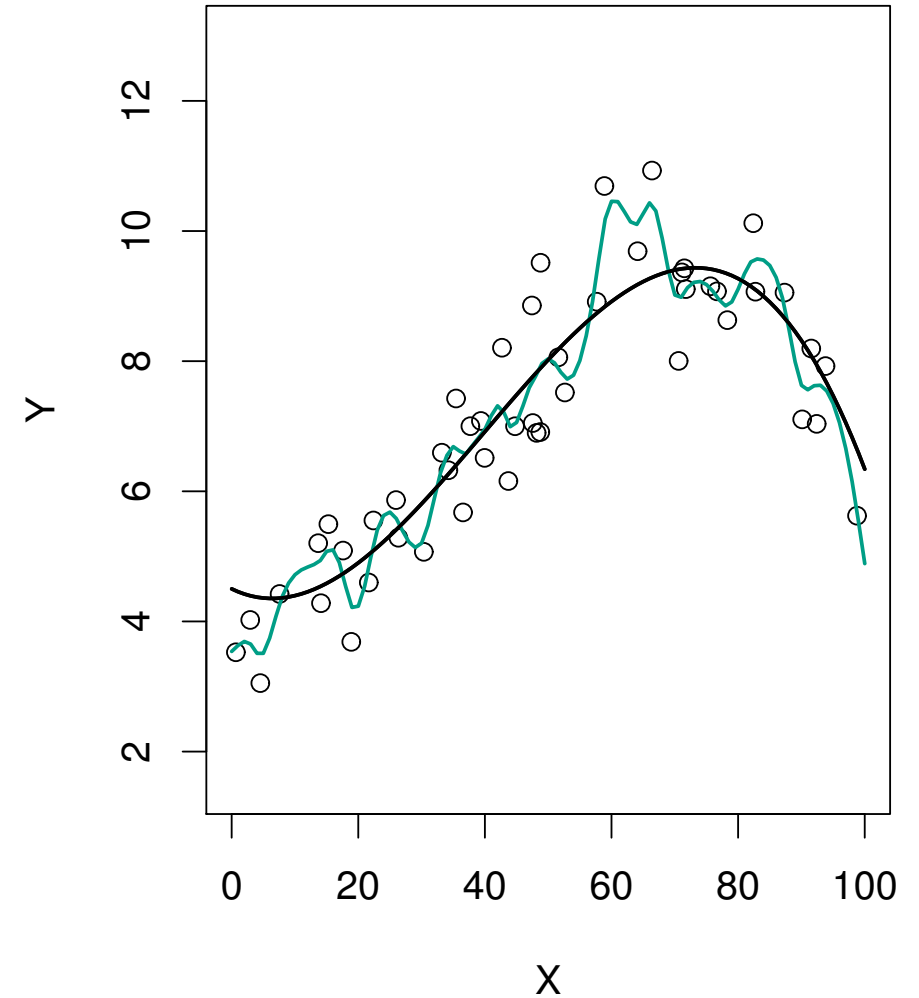
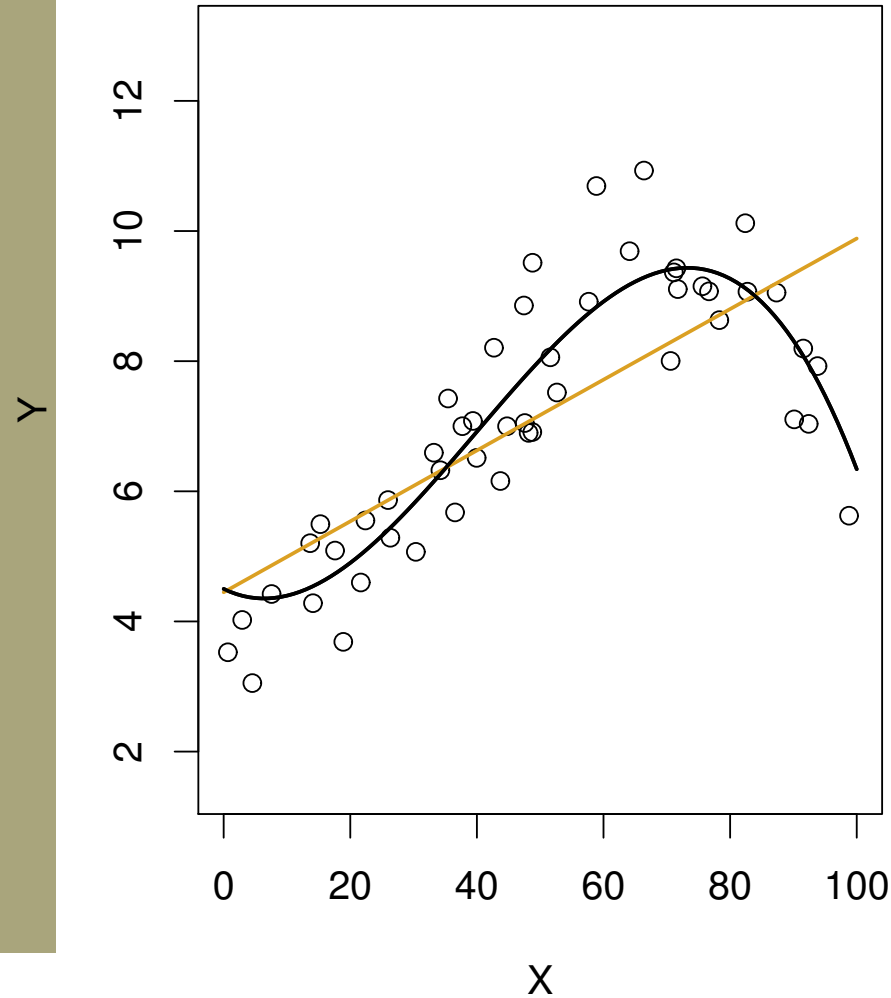
- In general, more flexible methods have **higher variance**



(**Variance**: the amount the model would change if we had different training data)

# Relationship between bias and variance

- In general, more flexible methods have **lower bias**



(**Bias:** the error introduced by approximating a complex phenomenon using a simple model)

## Trade-off between bias and variance

- Expected test MSE can be decomposed into three terms:

$$E\left(y_0 - \hat{f}(x_0)\right)^2 = \text{Var}\left(\hat{f}(x_0)\right) + \left[\text{Bias}\left(\hat{f}(x_0)\right)\right]^2 + \text{Var}(\varepsilon)$$

The variance of our model  
on the test value

The bias of our model  
on the test value

The variance  
of the error terms

# Balancing bias and variance

- It's easy to build a model with  
**low variance** but **high bias** (how?)
- Just as easy to build one with  
**low bias** but **high variance** (how?)
- The challenge: finding a method for which both the variance and the squared bias are low
- This trade-off is one of the most important recurring themes in this course

(**Variance**: the amount the model would change if we had different training data)

**Bias**: the error introduced by approximating a complex phenomenon using a simple model)

# What about classification?

- So far: how to evaluate a **regression** model
- Bias-variance trade-off also present in **classification**
- Need a way to deal with **qualitative responses**

What are some options?

# Training error rate

- **Common approach:** measure the proportion of the times our model incorrectly classifies a training data point

and take the average →  $\frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i)$

↑  
tally up all the times

where the model's classification was **different** from the true class

# Training error rate

- **Common approach:** measure the proportion of the times our model incorrectly classifies a training data point

and take the average  $\rightarrow \frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i)$

tally up  
all the times

where the model's  
classification was **different**  
from the true class

$$n = 3$$

| Student | Grade | C(X) |
|---------|-------|------|
| Ab      | B-    | C+   |
| Kaden   | B     | B    |
| Kylee   | A     | D    |

$$\begin{aligned} \text{Training error} &= \frac{1}{3} \sum_{i=1}^3 I(y_i \neq \hat{y}_i) \\ &= \frac{1}{3} ((1) + (0) + (1)) \\ &= \frac{2}{3} \\ &= 0.67 \end{aligned}$$

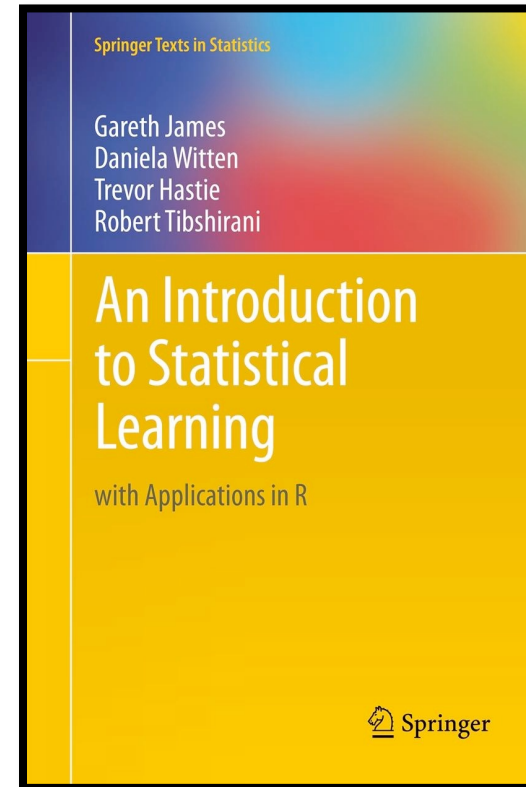


# Takeaways

- Choosing the “right” level of flexibility is **critical** (in both regression and classification)
- Bias-variance trade off makes this challenging
- Coming up in Ch. 5:
  - Various methods for **estimating** test error rates
  - How to use these estimates to find the **optimal level** of flexibility

# Reading

- In today's class, we covered ISLR: p. 29-37
- Next class, we'll have a crash course in linear regression (ISLR: p. p.59-82)

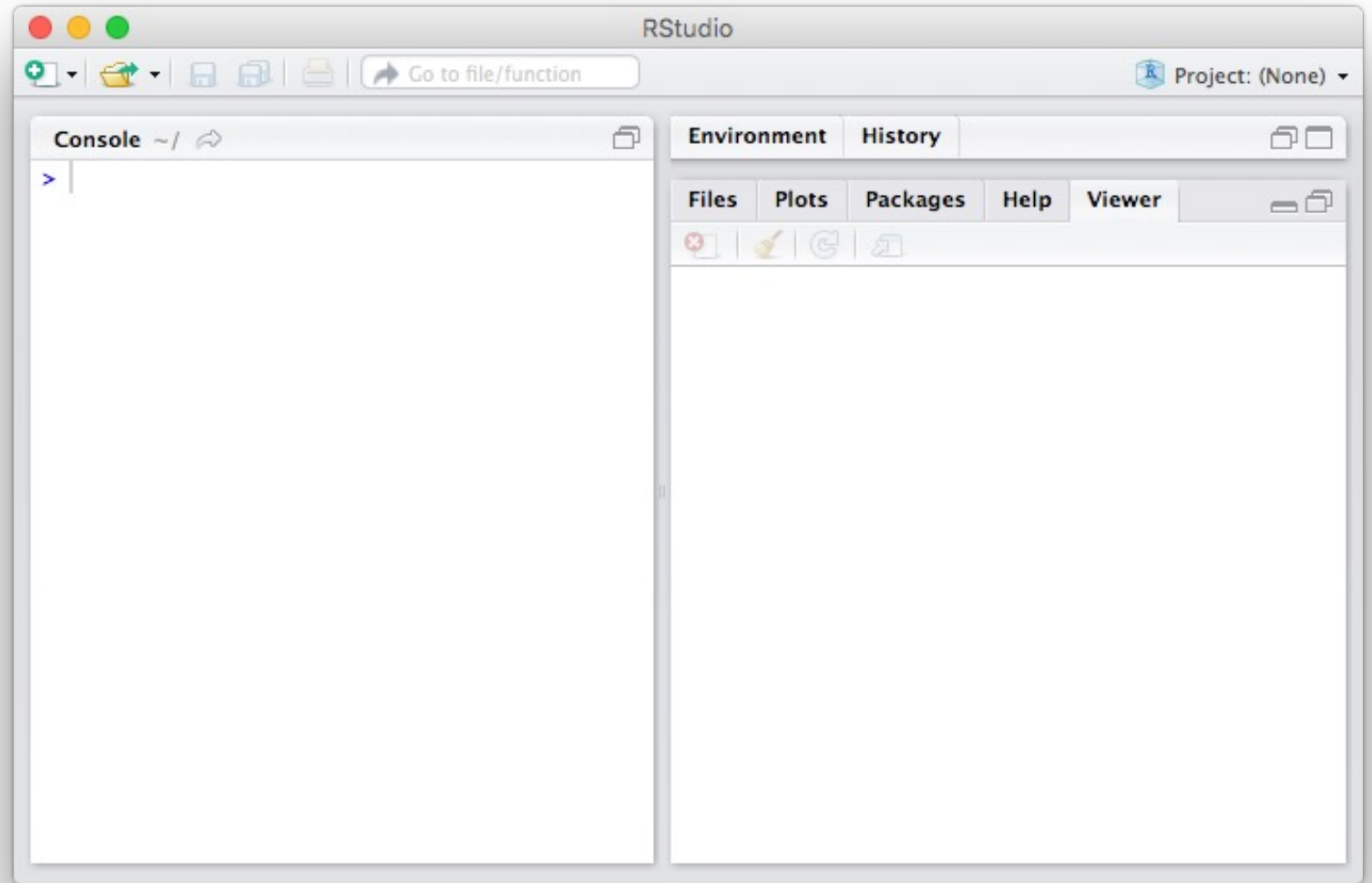


# Introduction to R



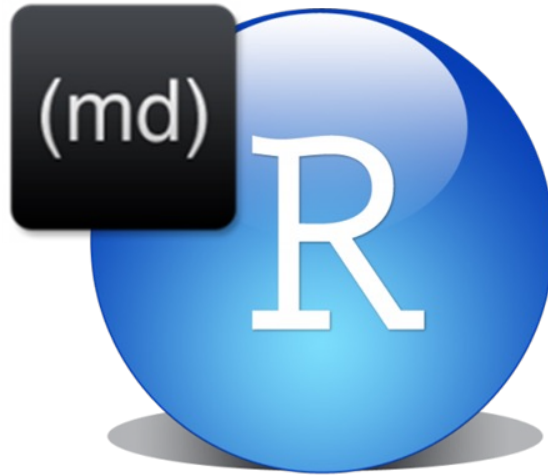
- Basic commands
- Loading external data
- Data wrangling 101
- Graphics
- Generating summaries

# Introduction to R



# Introduction to R

- Today's walkthrough was run using R Markdown:



- This allows me to build “notebooks” to combine step-by-step code and instructions/descriptions