

Math 300 Lesson 16 Notes
Multiple Regression - Related Topics

YOUR NAME HERE

June, 2022

Contents

Objectives 1
Reading 1
Lesson 1
Documenting software 6

Objectives

- 1. Find and interpret R-squared.
- 2. Use R-squared to select the best model.
- 3. Explain Simpson’s paradox and confounding variables.

Reading

Chapter 6.3

Lesson

There are no learning checks for this lesson, so work through the code and ideas in the reading.

- The book discusses Occam’s razor, which states that given the choice between a complex model and simple model pick the simple. This assumes similar performance. One way to measure performance is with R-squared. This is an internal measure of performance, it uses the data that was used to build the model. In a machine learning class, students can learn other ways to pick a model.
- R-squared always increases when we add a variable. Adjusted R-squared puts in a penalty for adding a variable. This penalty is just a heuristic. Analysts often use adjusted R-squared instead of R-squared.
- We preface our interpretation with the statement, “taking into account all the other explanatory variables in our model” in this section. This means we have to treat the other variables as at a constant value even though collinearity in practice may not allow this. It is only from an interpretation point of view that we use that statement.
- We must be aware of a phenomenon known as Simpson’s Paradox, whereby overall trends that exist in aggregate either disappear or reverse when the data are broken down into groups. We will discuss this today or next lesson.

Setup

```
library(tidyverse)
library(moderndiver)
library(skimr)
library(ISLR)
```

Recreate the analysis done in the book.

```
# Complete the code and remove comment symbols
# model_2_interaction <- lm(average_sat_math ~ _____ * size,
#                           data = MA_schools)
# get_regression_table(model_2_interaction)
```

```
# Complete the code and remove comment symbols
# model_2_parallel_slopes <- lm(average_sat_math ~ perc_disadvan _____ size,
#                               data = MA_schools)
# get_regression_table(model_2_parallel_slopes)
```

In a future lesson we will use the p-value and confidence intervals to determine the statistical importance of an explanatory variable.

```
# Complete the code and remove comment symbols
# get_regression_summaries(_____)
```

```
# Complete the code and remove comment symbols
# get_regression_summaries(_____)
```

R-squared

- Use R-squared to determine the model for the UT Austin teacher evaluation problem. (Objective 1 and 2)

Get the data

```
# Complete the code and remove comment symbols
# evals_ch6 <- _____ %>%
#   select(ID, score, age, gender)
```

Let's look at 5 random rows of data.

```
# Complete the code and remove comment symbols
# set.seed(941)
# evals_ch6 %>%
#   sample_n(size = _____)
```

- Interaction Model

In this model we allow a different slope and intercept for each gender.

```
# Complete the code and remove comment symbols
# ggplot(evals_ch6, aes(x = age, y = score, color = _____)) +
#   geom_point() +
#   labs(x = "Age", y = "Teaching Score", color = "Gender") +
#   geom_smooth(method = "lm", se = _____) +
#   theme_bw()
```

To get the model in R, we use the `*` which is not multiplication but an interaction term in the model formula.

```
# Complete the code and remove comment symbols
# Fit regression model:
# score_model_interaction <- lm(score ~ _____ * _____, data = evals_ch6)
```

```
# Complete the code and remove comment symbols
# Get regression table:
# get_regression_table(_____)
```

-
- Parallel Slopes Model

We will use the same data, but just build a different model.

```
# Complete the code and remove comment symbols
# ggplot(evals_ch6, aes(x = age, y = score, color = _____)) +
#   geom_point() +
#   labs(x = "Age", y = "Teaching Score", color = "Gender") +
#   geom_parallel_slopes(se = FALSE) +
#   theme_bw()
```

- Notice that the line for females stops at the extremes of the observed data.

```
# Complete the code and remove comment symbols
# Fit regression model:
# score_model_parallel_slopes <- lm(score ~ age _____ gender, data = evals_ch6)
```

```
# Complete the code and remove comment symbols
# Get regression table:
# get_regression_table(_____)
```

Now let's use R-squared to pick the model. We will use `rbind()` and `tidyverse` commands to put the results in a nice form.

```
# Complete the code and remove comment symbols
# get_regression_summaries(score_model_interaction) %>%
#   rbind(get_regression_summaries(_____)) %>%
#   mutate(model=c("Interaction", "Parallel Slopes")) %>%
#   select(model, r_squared, adj_r_squared)
```

Neither model is great since the R-squared is so small, we are only explaining 5 percent of the variation in the teacher score with our model. However, the more complex interaction model is better than the parallel slope model.

Simpson's paradox

It is key in modeling to account for lurking variable. This site Simpson's paradox.

Here is a nice example from the `palmerpenguins` data package.

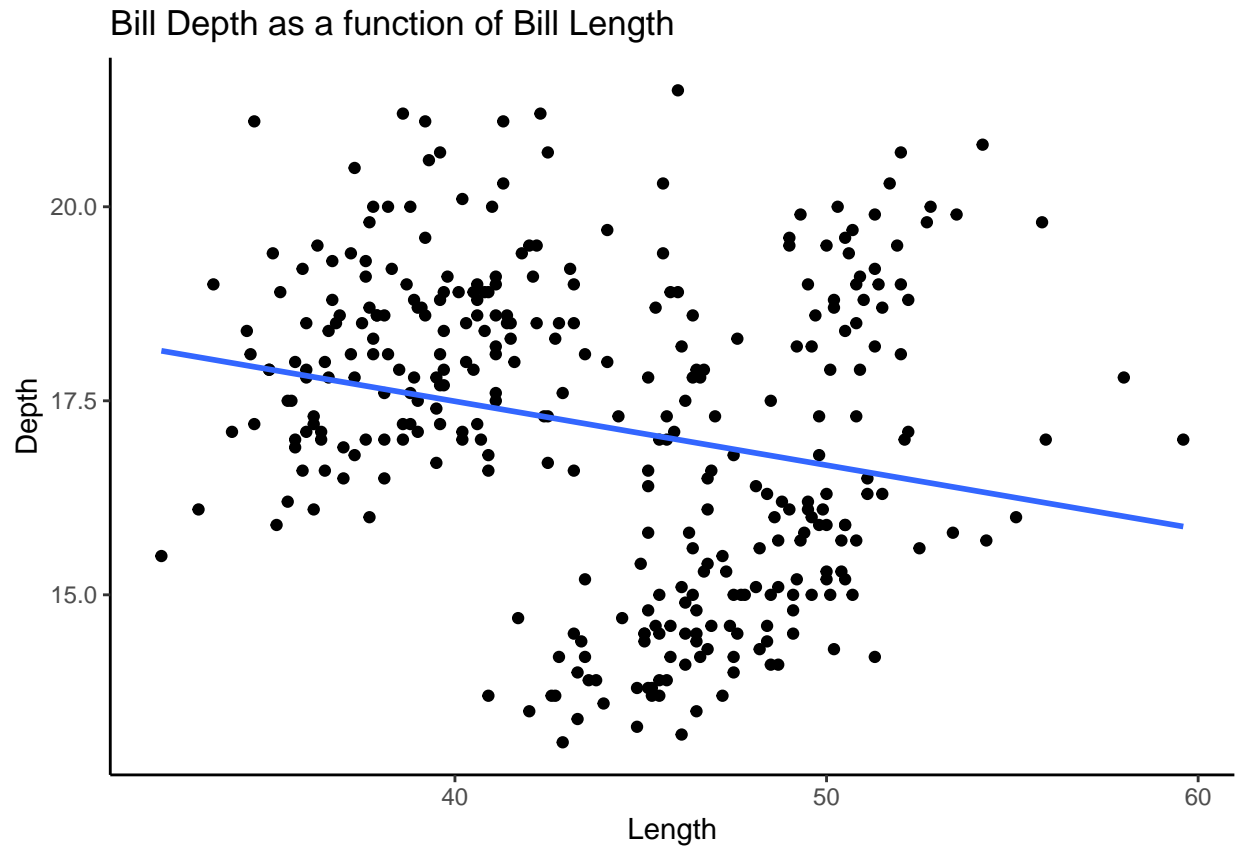
```
library(palmerpenguins)
```

```
penguin_df<-  
  palmerpenguins::penguins %>%  
  na.omit()
```

```
head(penguin_df)
```

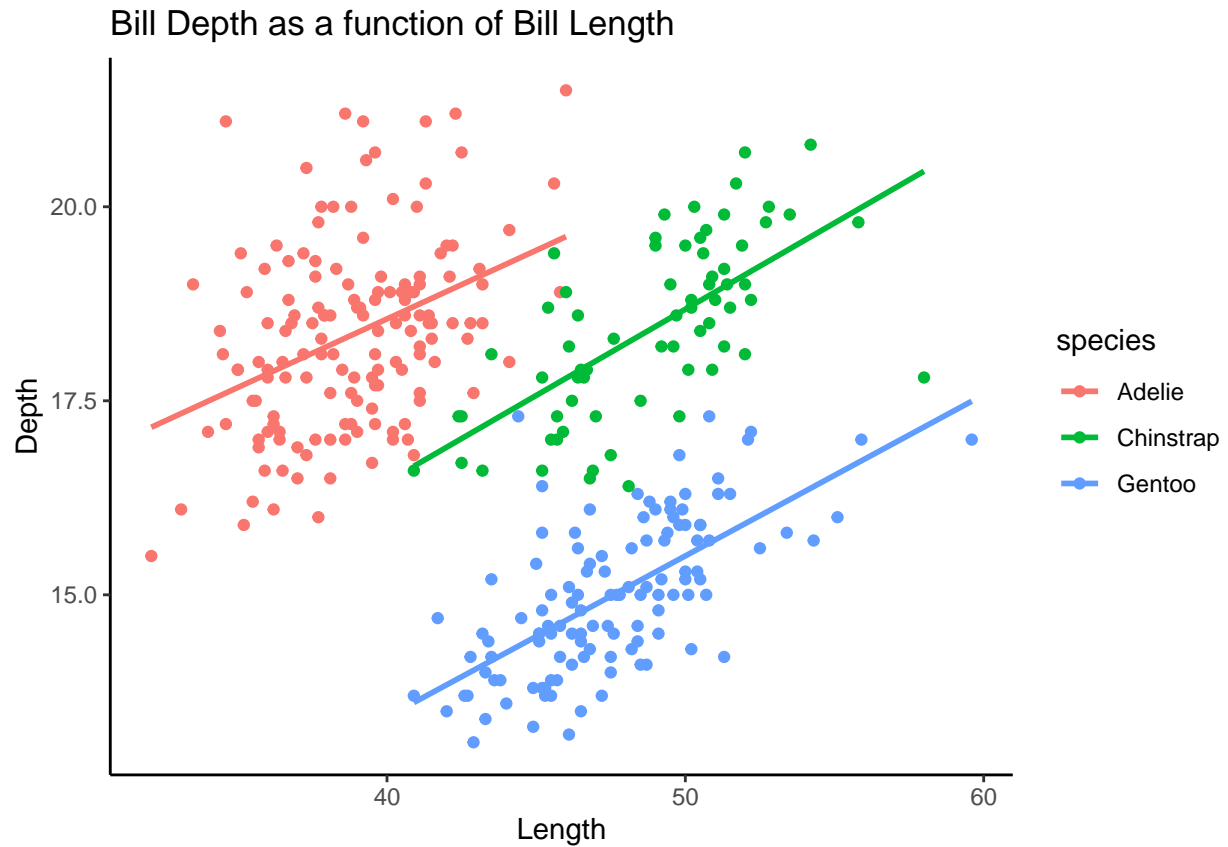
```
## # A tibble: 6 x 8  
##   species island bill_length_mm bill_depth_mm flipper_length_~ body_mass_g sex  
##   <fct>   <fct>         <dbl>         <dbl>         <int>         <int> <fct>  
## 1 Adelie  Torge~           39.1           18.7           181           3750 male  
## 2 Adelie  Torge~           39.5           17.4           186           3800 fema~  
## 3 Adelie  Torge~           40.3            18           195           3250 fema~  
## 4 Adelie  Torge~           36.7           19.3           193           3450 fema~  
## 5 Adelie  Torge~           39.3           20.6           190           3650 male  
## 6 Adelie  Torge~           38.9           17.8           181           3625 fema~  
## # ... with 1 more variable: year <int>
```

```
penguin_df %>%  
  ggplot(aes(x=bill_length_mm, y=bill_depth_mm)) +  
  geom_point() +  
  labs(x="Length", y="Depth", title="Bill Depth as a function of Bill Length") +  
  theme_classic() +  
  geom_smooth(method = "lm", se = FALSE)
```



From this we might as well conclude that the longer the bill, the less deep it is. However, if you drill down from the population level to the species level we see the opposite result.

```
penguin_df %>%  
  ggplot(aes(x=bill_length_mm, y=bill_depth_mm,  
             color=species)) +  
  geom_point() +  
  labs(x="Length", y="Depth", title="Bill Depth as a function of Bill Length") +  
  theme_classic() +  
  geom_smooth(method = "lm", se = FALSE)
```



Explain these results in terms of a confounding variable.

Documenting software

- File creation date: 2022-06-27
- R version 4.1.3 (2022-03-10)
- tidyverse package version: 1.3.1
- skimr package version: 2.1.4
- palmerpenguins package version: 0.1.0
- moderndive package version: 0.5.4