

# Final Project

Your Name Here

2023-11-28

## Final Details

### Purpose

The goal of this final assignment is to assess your ability to integrate the many skills you have learned over the semester: filtering and summarizing data, creating new columns, choosing the appropriate data visualizations, and performing and interpreting the appropriate statistical tests.

### Task

Write R code which produces the correct data, summaries, plots and analyses. Correctly interpret the results of these plots and analyses.

### Criteria for Success

- Code chunks run without errors
- Code produces the correct result
  - Code that produces the correct answer will receive **full** credit
  - Code attempts with logical direction will receive **partial** credit
- Appropriate plot types are used to visualize the data
- Appropriate statistical tests are used to analyze the data
- Written answers address the questions in sufficient detail

### Due Date

Dec 11 at midnight MST

## Final

For your final this semester, I am presenting you with 3 problem sets, totaling up to 60 points.

I'm expecting you to be able to filter and summarize the data in ways you need, choose the appropriate visualization, choose the appropriate analysis, and correctly interpret the analysis for the question I've asked you.

It is important to note that I will not be giving you an answer key for this final project since a major part of your grade for the assignment is being able to choose the appropriate visualizations and analyses for the question and the data.

We are going to use the `palmerpenguins` R package and data set, which we've used many times before! You can learn more about it [here](#). This is a real data set from a Long-Term Ecological Research (LTER) site in Antarctica.

---

## Set-Up

**Be sure to run both of these code chunks before you begin!** I've gone ahead and included the code to load the two packages you will need to successfully complete this project: the `tidyverse` and `palmerpenguins`. Be sure to run this code chunk!

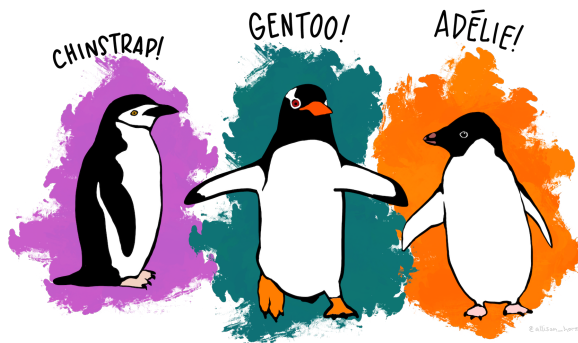
```
library(tidyverse)
library(palmerpenguins)
```

**Important!** I've also included one more code chunk below. Be sure to run this code chunk, as well! It does two key things:

- first, it adds the `penguins` data frame to your environment, which I imagine you will find helpful
- second, it removes all rows that have any NA values, which will make completing this assignment a bit easier.

Once you've run this line of code, you should see the `penguins` data frame pop up in your environment with 333 observations (rows) and 8 variables (columns).

```
penguins <- penguins %>% drop_na()
```



---

## Structure & Guidelines

Like the practice version, this final project is structured as 3 different problem sets. For each problem set, I am presenting you with an initial question to guide your thinking and analysis.

## Data

Assume that nothing carries over between problem sets.

Each problem set is stand-alone, meaning that you should always start with the `penguins` data frame at the beginning of each problem set. If you should use a data frame that you created *within* the problem set, I explicitly state so.

For example, in Problem Set 2, you should use the `biscoe` data frame that you create for the entire problem set; at the start of Problem Set 3, start over with the `penguins` data frame.

## Interpreting Statistical Results

When I ask you to interpret statistical results, you should roughly follow these guidelines.

- the cut-off for our p-values is always 0.05
- report the p-value that we are focused on
- if there are multiple p-values of interest, report all of them
- state whether the p-value indicates a significant difference/relationship
- if applicable, state whether we should or should not reject the null hypothesis

## Plotting

All plots should be made using `ggplot2`.

Your options for plot types to choose from are:

- multiple histogram plots
  - use transparency (`alpha`)
  - use `position = "identity"` with multiple groups to see the full distributions
- multiple density plots
  - use transparency (`alpha`) with multiple groups to see the full distributions
- box-and-whisker plot
  - add points on top of the box plot to show the distribution of the points
- scatter plot
  - add the linear model to every scatter plot

**Note: All plots should have modified axis labels and legend labels.**

In many cases, this might mean capitalizing the axis label or legend label. In other cases, you might want to put units in parentheses after the words (e.g., Body Mass (g)).

## Problem Set 1 (15 points)

**Question:** Are there differences in the average bill length across the 3 islands in the data set: Dream, Biscoe, and Torgersen? (Ignore species) Let's start by summarizing the bill length data.

1. Calculate at least one measure of central tendency and one measure of variability for the bill length for *each* island. (2 points)

```
# measures of central tendency include mean, median, or mode
# measures of dispersion are st dev or BOTH min and max
penguins %>%
  group_by(island) %>%
  summarise(min_bill = min(bill_length_mm),
            max_bill = max(bill_length_mm),
            mean_bill = mean(bill_length_mm))
```

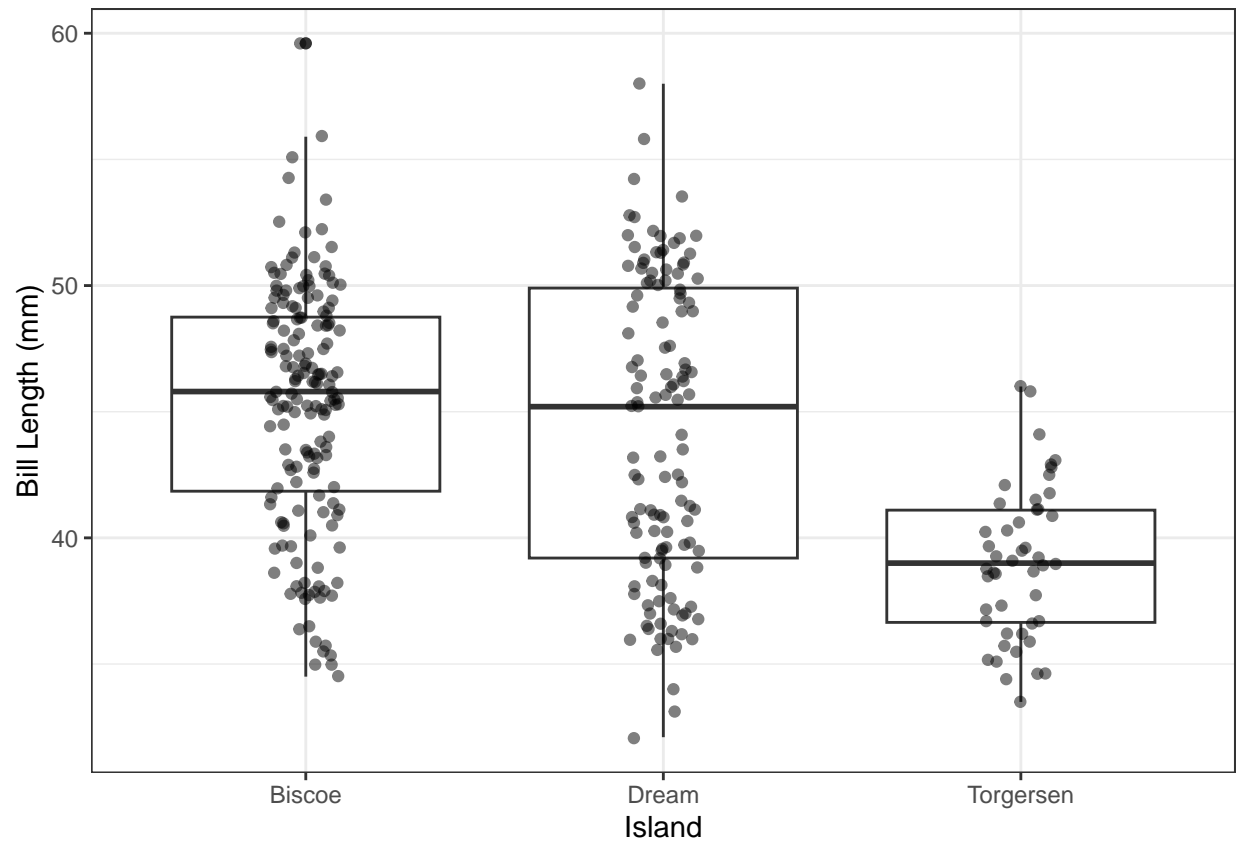
```
## # A tibble: 3 x 4
##   island   min_bill max_bill mean_bill
##   <fct>     <dbl>   <dbl>     <dbl>
## 1 Biscoe    34.5     59.6      45.2
## 2 Dream    32.1     58       44.2
## 3 Torgersen 33.5     46       39.0
```

2. Which of our variables would be considered *independent* and which one *dependent*? Also determine whether each is *continuous* or *categorical*. (2 points)
  - **island**: independent, categorical
  - **bill length**: dependent, continuous

Now that we have an idea numerically of the differences between the islands, let's plot the differences.

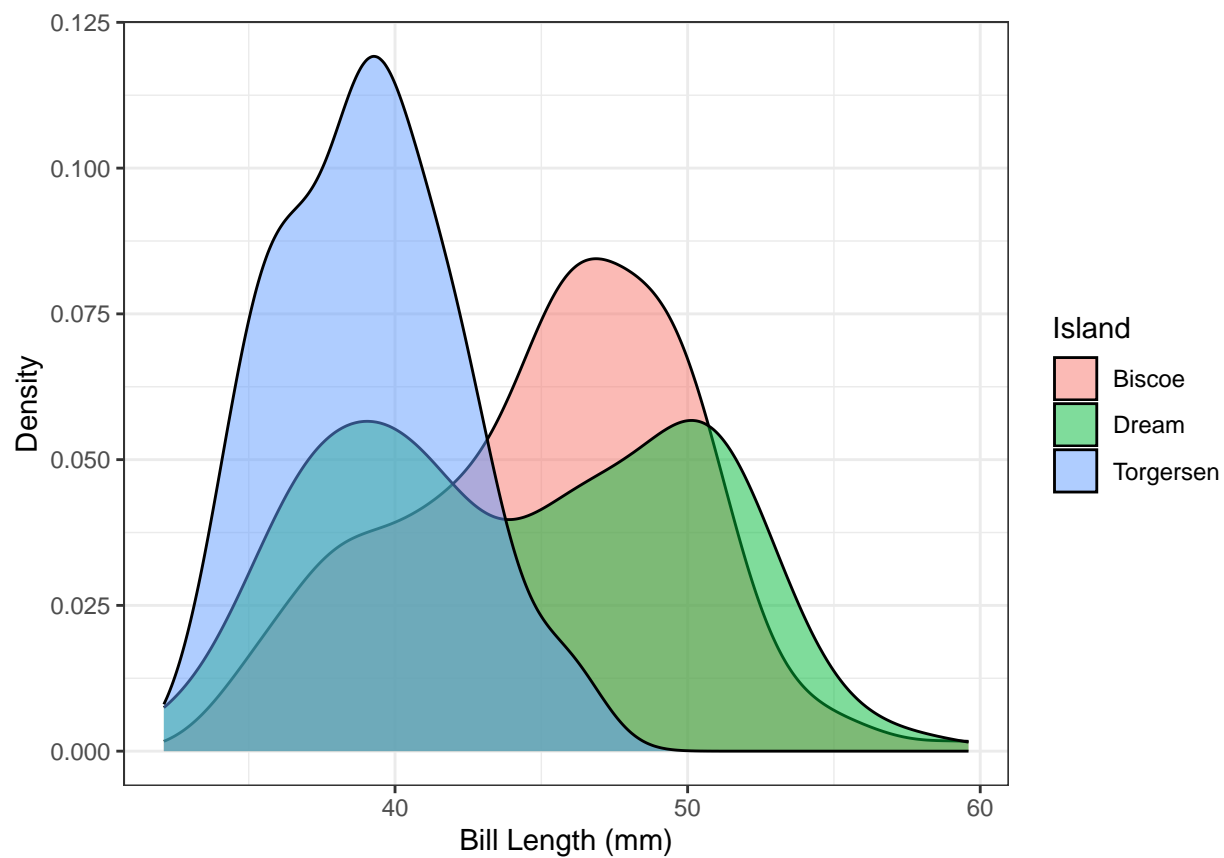
3. Choose an appropriate plot. Ensure that you follow the plotting guidelines in the Structure & Guidelines section above! (2 points)

```
ggplot(penguins, aes(island, bill_length_mm)) + # adding color is optional
  geom_boxplot() +
  geom_jitter(width = 0.1, alpha = 0.5) +
  labs(x = "Island",
       y = "Bill Length (mm)") +
  theme_bw()                                     # if adding color, must change color label
                                                # can choose any theme but must include
```



```
# OR

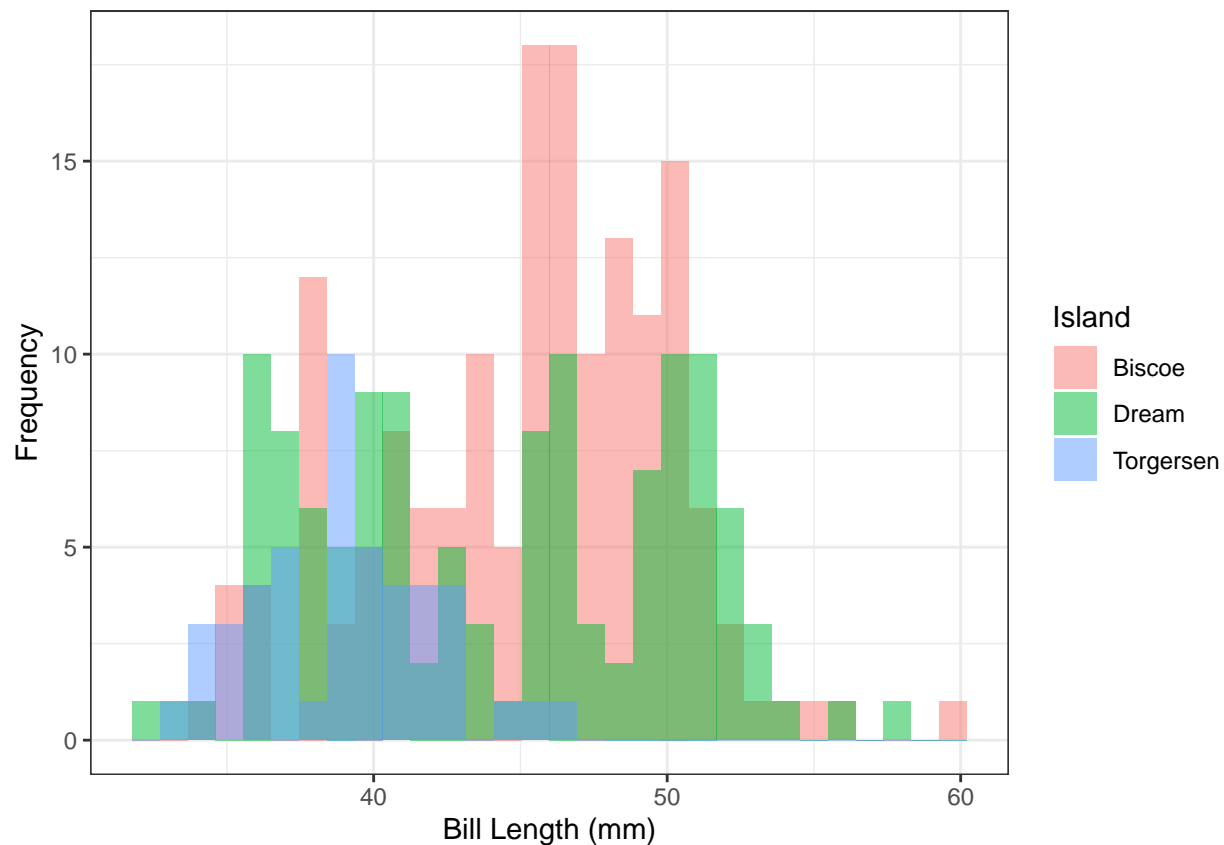
ggplot(penguins, aes(bill_length_mm, fill = island)) +
  geom_density(alpha = 0.5) +
  labs(x = "Bill Length (mm)",
       y = "Density",
       fill = "Island") +
  theme_bw()
```



*# OR*

```
ggplot(penguins, aes(bill_length_mm, fill = island)) +  
  geom_histogram(alpha = 0.5, position = "identity") +  
  labs(x = "Bill Length (mm)",  
       y = "Frequency",  
       fill = "Island") +  
  theme_bw()
```

## 'stat\_bin()' using 'bins = 30'. Pick better value with 'binwidth'.



4. Write the correct pair of statistical hypotheses for our question. (2 points)

**Null:** no difference in the mean bill length between penguins on different islands

**Alternative:** true difference in the mean bill length between penguins on different islands

5. Run the appropriate statistical analysis for our question. (2 points)

```
aov_model <- aov(data = penguins, bill_length_mm ~ island)
summary(aov_model)
```

```
##           Df Sum Sq Mean Sq F value    Pr(>F)
## island      2   1417    708.6   27.47 9.21e-12 ***
## Residuals 330    8512     25.8
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

6. Interpret the results of this test. (3 points)

- is there a significant difference?
- what does that significant difference mean?
- should we reject the null hypothesis?

*Answer: yes,  $p = 9.21e-12$ , which is smaller than 0.05, reject null. penguins on different islands have significantly different beak lengths*

7. Should we run pairwise comparisons? If no, explain why not. If yes, do so and interpret the results. (2 points)

*Answer:* yes; significant differences between all pairs except Dream-Biscoe because  $p > 0.05$

```
TukeyHSD(aov_model)
```

```
## Tukey multiple comparisons of means
## 95% family-wise confidence level
##
## Fit: aov(formula = bill_length_mm ~ island, data = penguins)
##
## $island
##              diff      lwr      upr    p adj
## Dream-Biscoe -1.026515 -2.454611  0.4015809 0.2096455
## Torgersen-Biscoe -6.210168 -8.189815 -4.2305220 0.0000000
## Torgersen-Dream -5.183653 -7.234077 -3.1332293 0.0000000
```

---

## Problem Set 2 (30 points total)

This problem set has 2 parts.

### Part 1 (20 points)

**Question: Is there a significant relationship between bill length and bill depth for penguins on Biscoe Island?** For this problem set, we are going to use data from Biscoe island only.

1. Create a new data frame called `biscoe` that includes only penguins from Biscoe island. This new data frame should have 163 rows. (2 points)

```
biscoe <- penguins %>%
  filter(island == "Biscoe")
```

You will want to use the `biscoe` data set for the rest of Problem Set #2.

This is a scenario where there is no independent and no dependent variable. Go ahead and **treat bill length as the independent variable (x-axis) and bill depth as the dependent variable (y-axis).**

2. Determine whether bill length and bill depth are continuous or categorical (2 points)
  - **bill depth:** dependent,
  - **bill length:** independent

For now, ignore species. We will address species in Part 2 of the problem set.

3. Write the correct pair of statistical hypotheses for our question. (2 points)

**Null:** no relationship between bill length and bill depth

**Alternative:** there is a relationship between bill length and bill depth

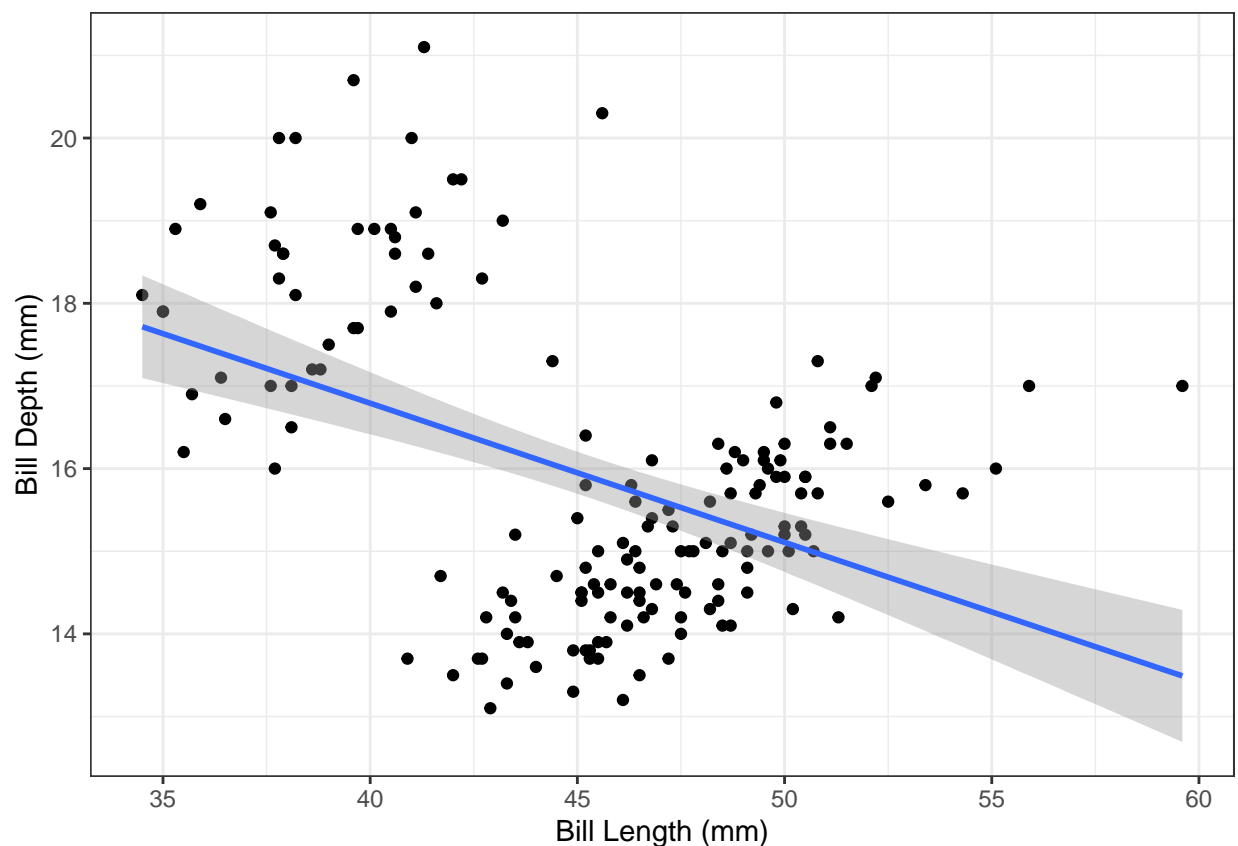
4. Plot the relationship between bill length and bill depth using the appropriate plot type. (3 points)



- Be sure to add a line of best fit using the `geom_smooth` function—and make sure it is a straight line (no wiggles, which the default will produce).
- Ensure that the plot has clear labels on the axes (follow the Structure & Guidelines).
- Remember, we are ignoring species for now.

```
ggplot(biscoe, aes(bill_length_mm, bill_depth_mm)) +
  geom_point() +
  geom_smooth(method = "lm") +
  labs(x = "Bill Length (mm)",
       y = "Bill Depth (mm)") +
  theme_bw()
```

```
## 'geom_smooth()' using formula = 'y ~ x'
```



5. Describe the relationship that you see in the plot you just made. Is the relationship positive, negative, or nonexistent? How do you know? (2 points)

*Answer:* negative; downward slope, high values associated with low values on the other axis

6. Calculate the correlation coefficient,  $r$ . What does this value confirm for us about the relationship (positive, negative, no relationship)? (2 points)

```
r <- cor(biscoe$bill_length_mm, biscoe$bill_depth_mm)
r
```

```
## [1] -0.4446658
```

Answer: negative relationship

7. Calculate the  $r^2$  value. How much variation is explained by the line of best fit? Remember, this number is typically expressed as a percent (x 100). (2 points)

```
r^2 * 100
```

```
## [1] 19.77277
```

Answer: About 20% variation explained

8. Let's see if there is a significant relationship between bill length and bill depth. Perform the correct statistical analysis (1 point) and interpret the results. (5 points total)

- What is the equation of the line of best fit? Use both variable names and values from the statistical analysis. (1 point)
- What is the p-value? (1 point)
- Is there a significant relationship? (1 point)
- What should we do with the null hypothesis? (1 point)

```
summary(lm(data = biscoe, bill_depth_mm ~ bill_length_mm))
```

```
##
## Call:
## lm(formula = bill_depth_mm ~ bill_length_mm, data = biscoe)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.2027 -1.2536 -0.1135  1.0735  4.5279
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   23.52509     1.21614   19.344 < 2e-16 ***
## bill_length_mm -0.16835     0.02673   -6.299 2.74e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.642 on 161 degrees of freedom
## Multiple R-squared:  0.1977, Adjusted R-squared:  0.1927
## F-statistic: 39.68 on 1 and 161 DF, p-value: 2.738e-09
```

Answer:  $bill\_depth = -0.168 * bill\_length + 23.5$ ;  $p = 2.74e-9$ , highly significant

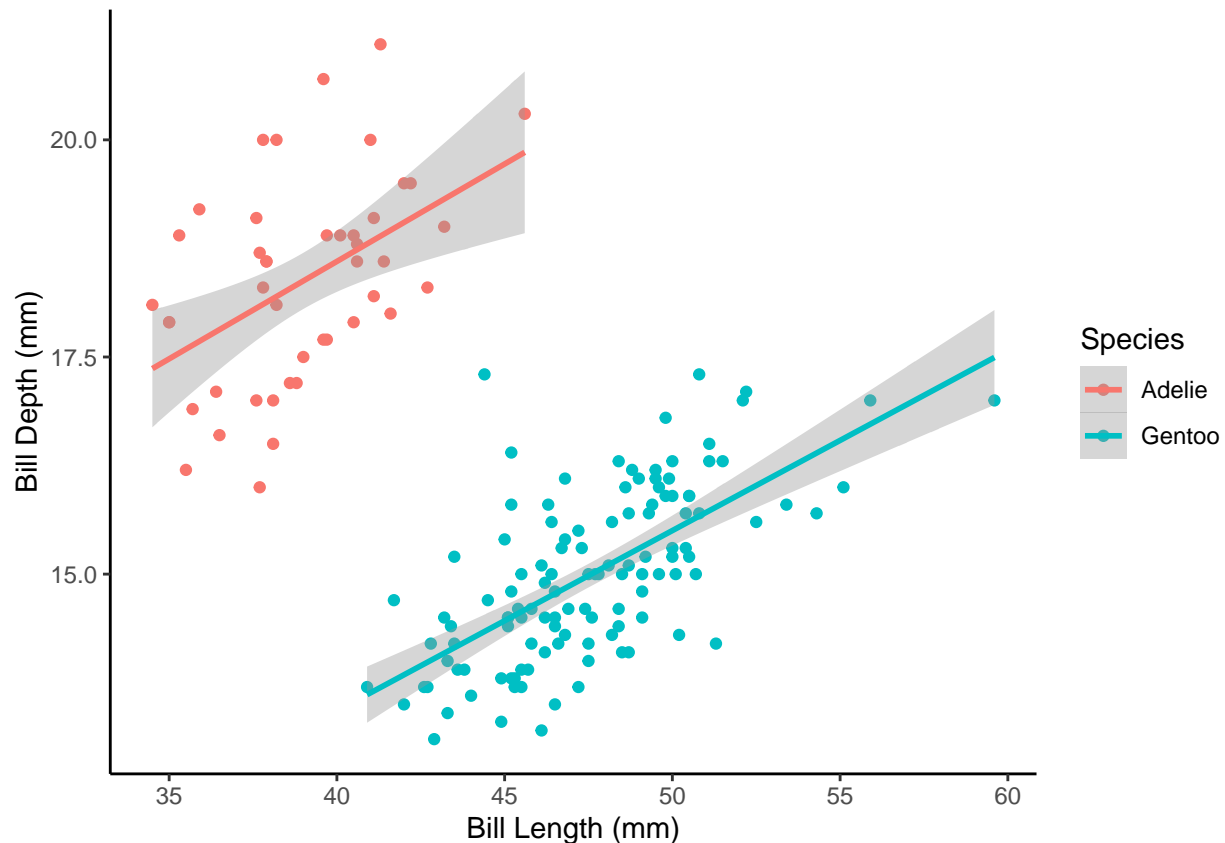
## Part 2 (10 points)

When we look at the plot of the data, it looks like there might be two different groups in the data. Let's see what happens when we add in species to this analysis.

9. Let's make the color of the points and the linear models differ by species on Biscoe Island. Be sure to adjust *all* labels on the plot accordingly. (3 points)

```
ggplot(biscoe, aes(bill_length_mm, bill_depth_mm, color = species)) +
  geom_point() +
  geom_smooth(method = "lm") +
  labs(x = "Bill Length (mm)",
       y = "Bill Depth (mm)",
       color = "Species") +
  theme_classic()
```

```
## 'geom_smooth()' using formula = 'y ~ x'
```



10. Run the appropriate statistical test, adding species (and the interaction between the two independent variables) into the model. (2 points)

```
summary(lm(data = biscoe, bill_depth_mm ~ bill_length_mm * species))
```

```
##
## Call:
## lm(formula = bill_depth_mm ~ bill_length_mm * species, data = biscoe)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.08441 -0.64479 -0.02957  0.46954  2.96109
##
```

```
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      9.62634    2.02455   4.755 4.42e-06 ***
## bill_length_mm    0.22435    0.05184   4.328 2.66e-05 ***
## speciesGentoo    -4.50539    2.34915  -1.918  0.0569 .
## bill_length_mm:speciesGentoo -0.01674    0.05755  -0.291  0.7715
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.8434 on 159 degrees of freedom
## Multiple R-squared:  0.791, Adjusted R-squared:  0.7871
## F-statistic: 200.6 on 3 and 159 DF, p-value: < 2.2e-16
```

11. Interpret the results of the test above. Is species significant? Is there a significant interaction between bill length and species? (2 points)

*Answer: neither species nor interaction is significant*

12. Write 2-3 sentences discussing if and/or how adding species into our linear model changes our interpretation of the data. Did the type of relationship change? Did significance levels change? Do the two linear models tell us different things? (3 points)

*Answer: negative to positive; models still tell us similar things in terms of significance, though*

---

### Problem Set 3 (15 points)

**Question: Is there a difference in the average flipper length between male and female Chinstrap penguins?**

1. Our first step is to create a new data frame that includes only Chinstrap penguins. Call this new data frame `chinstrap`. (1 point)

```
chinstrap <- penguins %>%
  filter(species == "Chinstrap")
```

We will be using the `chinstrap` data frame for the rest of this problem set.

2. Let's summarize our data. Calculate one measure of central tendency and one (complete) measure of variability of the flipper length column for *each* sex: male and female. Save this dataframe as `chinstrap_summary`. (2 points)

```
chinstrap_summary <- chinstrap %>%
  group_by(sex) %>%
  summarise(mean_flipper = mean(flipper_length_mm),
            sd_flipper = sd(flipper_length_mm))
```

3. What if we wanted our summary data in centimeters instead of millimeters?
  - a. First, create a function that will convert a number from millimeters to centimeters. (1 point)

```
mm_to_cm <- function(mm) {
  cm <- mm / 10
  return(cm)
}
```

- b. Now, using the same code from question 3 above, but add one line (in the correct location) that uses your newly created `mm_to_cm` function and produces the same data frame but with the summary values in cm instead of mm. Made sure to edit the summary functions accordingly, as well. (1 point)

```
chinstrap_summary <- chinstrap %>%
  mutate(flipper_length_cm = mm_to_cm(flipper_length_mm)) %>%
  group_by(sex) %>%
  summarise(mean_flipper = mean(flipper_length_cm),
            sd_flipper = sd(flipper_length_cm))
```

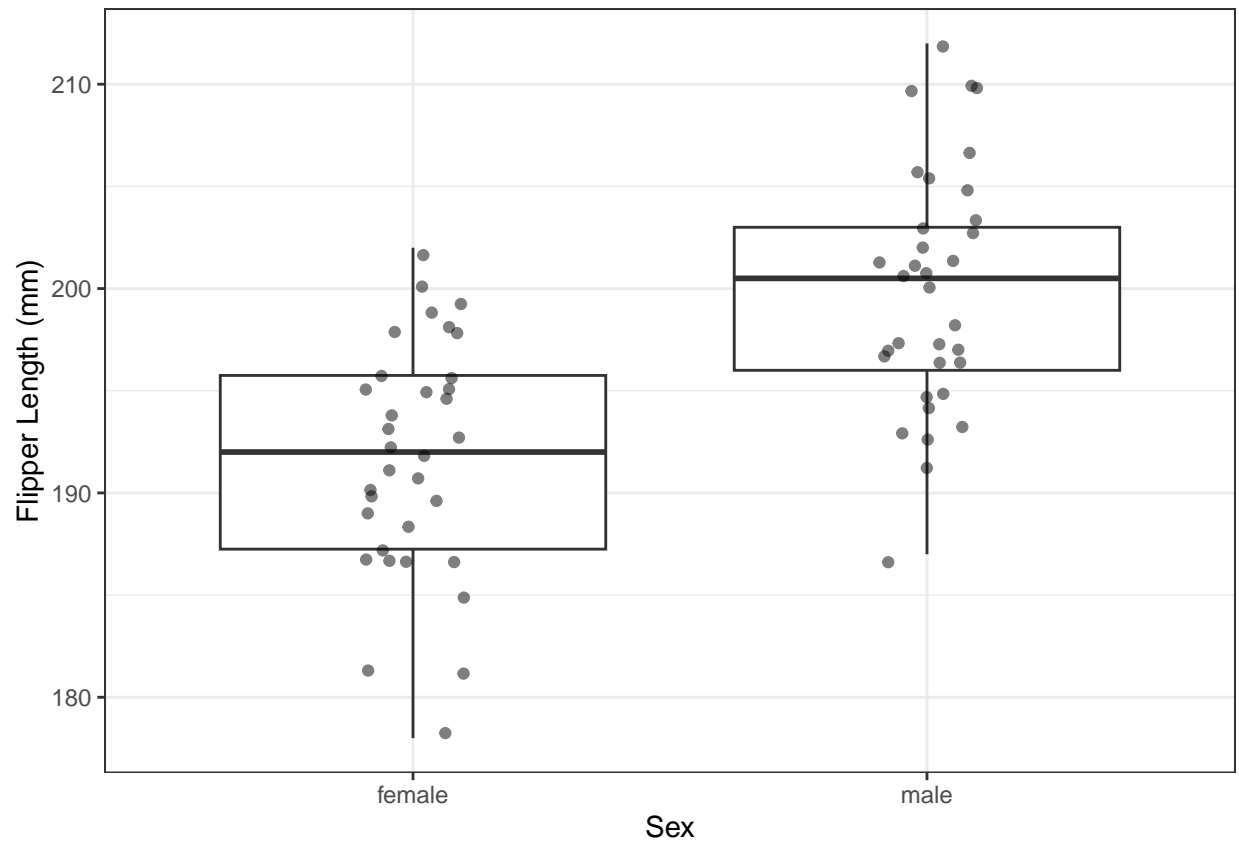
Ok, we have summarized the flipper length data for our two groups! Let's get back to the `chinstrap` dataframe (not the summary data frame), and keep working.

4. Determine which variable is dependent and which is independent. Also determine if each variable is continuous or categorical. (1 points)
  - **flipper length**: dependent, continuous
  - **species**: independent, categorical

Let's plot the body mass data for the two groups.

5. Choose an appropriate plot for data with one continuous variable and one categorical variable (there are a few options). Be sure to adjust the x- and y-axis labels appropriately. (2 points)

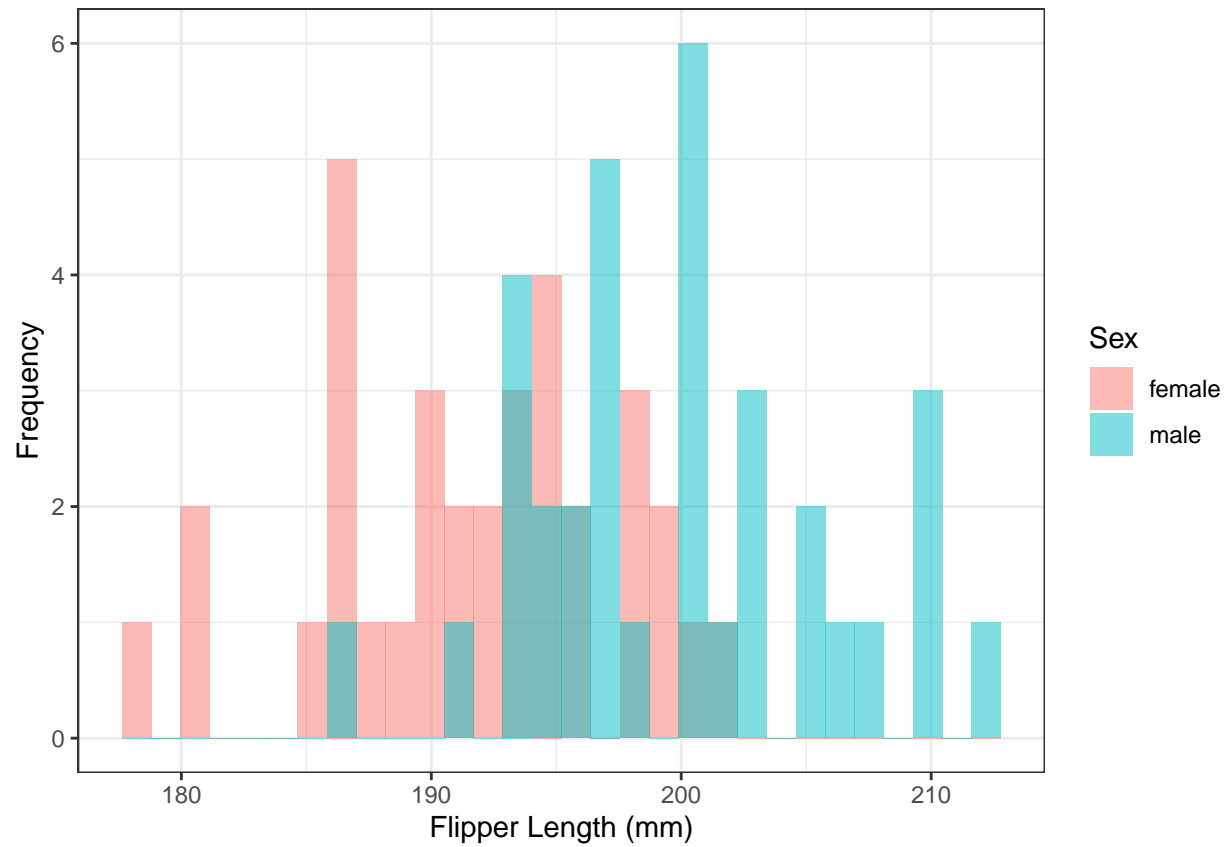
```
ggplot(chinstrap, aes(sex, flipper_length_mm)) +
  geom_boxplot() +
  geom_jitter(width = 0.1, alpha = 0.5) +
  labs(x = "Sex",
       y = "Flipper Length (mm)") +
  theme_bw()
```



*# OR*

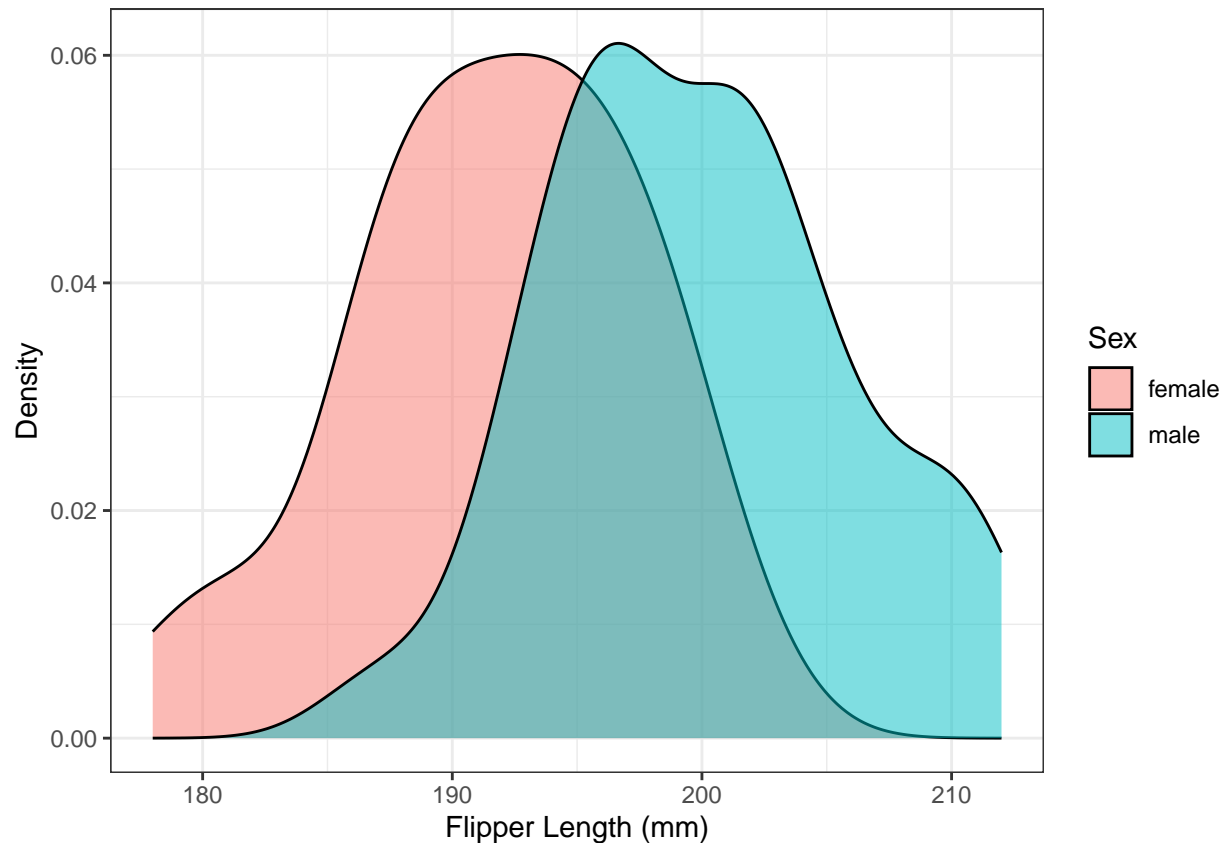
```
ggplot(chinstrap, aes(flipper_length_mm, fill = sex)) +
  geom_histogram(alpha = 0.5, position = "identity") +
  labs(x = "Flipper Length (mm)",
       y = "Frequency",
       fill = "Sex") +
  theme_bw()
```

## 'stat\_bin()' using 'bins = 30'. Pick better value with 'binwidth'.



```
# OR

ggplot(chinstrap, aes(flipper_length_mm, fill = sex)) +
  geom_density(alpha = 0.5) +
  labs(x = "Flipper Length (mm)",
       y = "Density",
       fill = "Sex") +
  theme_bw()
```



6. Write the pair of statistical hypotheses for our question. (1 points)

**Null:** no difference in the mean flipper length between male and female chinstrap penguins **Alternative:** true difference in the mean flipper length between male and female chinstrap penguins

7. Perform the appropriate analysis to compare the flipper lengths of each species. (2 points)

```
t.test(data = chinstrap, flipper_length_mm ~ sex)
```

```
##
##  Welch Two Sample t-test
##
## data:  flipper_length_mm by sex
## t = -5.7467, df = 65.905, p-value = 2.535e-07
## alternative hypothesis: true difference in means between group female and group male is not equal to
## 95 percent confidence interval:
##  -11.017272  -5.335669
## sample estimates:
## mean in group female    mean in group male
##           191.7353           199.9118
```

8. Interpret the results of this test. (2 points)

- is there a significant difference?
- what does that significant difference mean?
- should we reject the null hypothesis?



*Answer: yes, there is a significant difference ( $p = 2.535 \times 10^{-7}$ ); reject null*

9. Should we run pairwise comparisons? If no, explain why not. If yes, do so and interpret the results. (2 points)

*Answer: no, t-test is comparing just one pair*

## **The End!**

Great work, and thanks for a wonderful semester!