

## module\_3\_2

Keaton Wilson

2/12/2020

### Plotting and Means Comparisons

Last class, we examined the collars data set, read a bit about the grammar of graphics and analyzed some code for a ggplot of some of the data we've been working with.

Let's take a stab at making our ggplot from scratch.

In your groups, work together to develop code that utilizes the collars data <https://tinyurl.com/sp7b25x> and plots:

1. Two histograms one for battery life and one for signal length
2. Color coded for each manufacturer
3. Correct and appropriate labels
4. A vertical line that plots the median of each group (this is trickier - remember to use the google)

```
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.1 --
```

```
## v ggplot2 3.3.5      v purrr   0.3.4
## v tibble  3.1.6      v dplyr   1.0.7
## v tidyr   1.1.4      v stringr 1.4.0
## v readr   2.1.1      v forcats 0.5.1
```

```
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
```

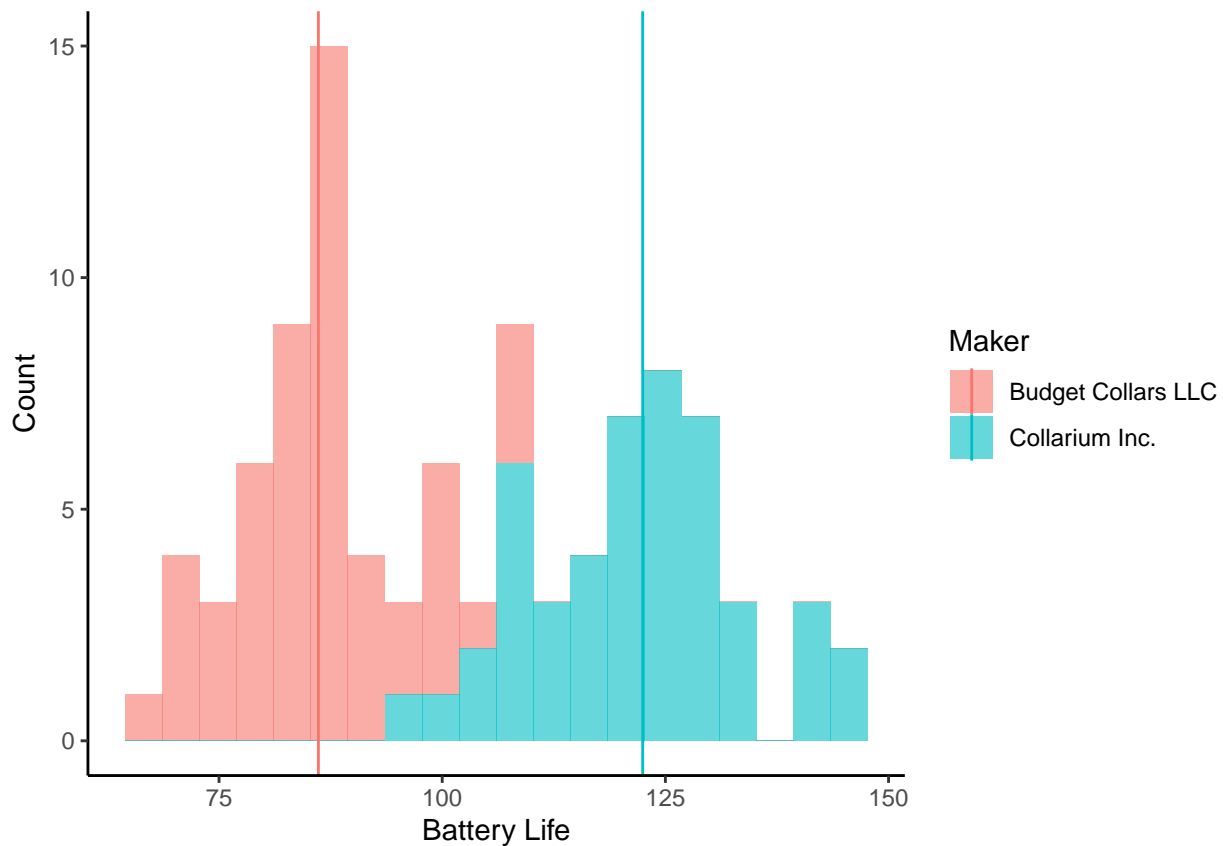
```
collars = read_csv("../data/collar_data.csv")
```

```
## Rows: 100 Columns: 5
```

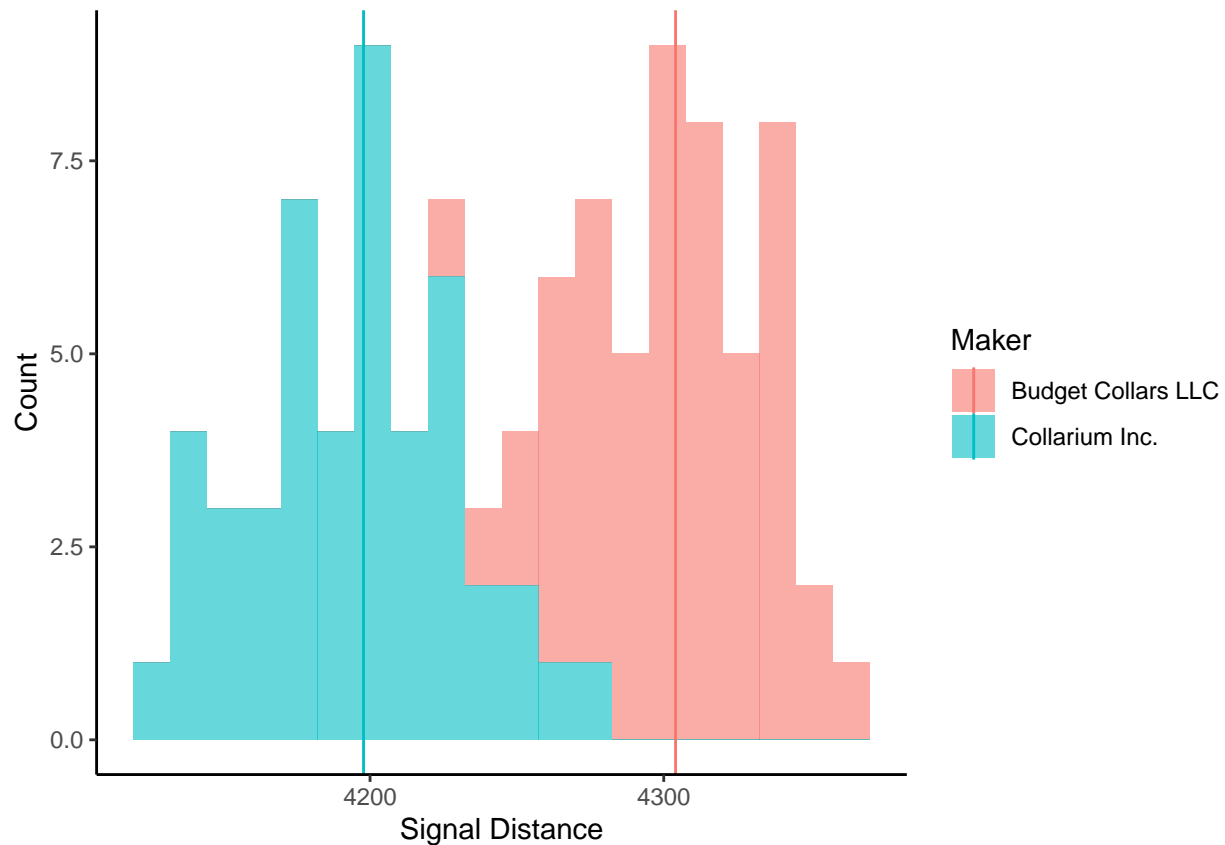
```
## -- Column specification -----
## Delimiter: ","
## chr (1): maker
## dbl (4): collar_id, battery_life, signal_distance, fail
```

```
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
```

```
collars %>%
  ggplot(aes(x=battery_life, fill = maker)) +
  geom_histogram(alpha = 0.6, bins = 20) +
  geom_vline(data = collars %>%
    group_by(maker) %>%
    summarize(median_battery = median(battery_life)),
    aes(xintercept = median_battery, col = maker)) +
  ylab("Count") +
  xlab("Battery Life") +
  scale_color_discrete(name = "Maker") +
  scale_fill_discrete(name = "Maker") +
  theme_classic()
```



```
collars %>%
  ggplot(aes(x=signal_distance, fill = maker)) +
  geom_histogram(alpha = 0.6, bins = 20) +
  geom_vline(data = collars %>%
    group_by(maker) %>%
    summarize(median_signal = median(signal_distance)),
    aes(xintercept = median_signal, col = maker)) +
  ylab("Count") +
  xlab("Signal Distance") +
  scale_color_discrete(name = "Maker") +
  scale_fill_discrete(name = "Maker") +
  theme_classic()
```



What is another way to plot these data?

## Stepping back and thinking about the problem

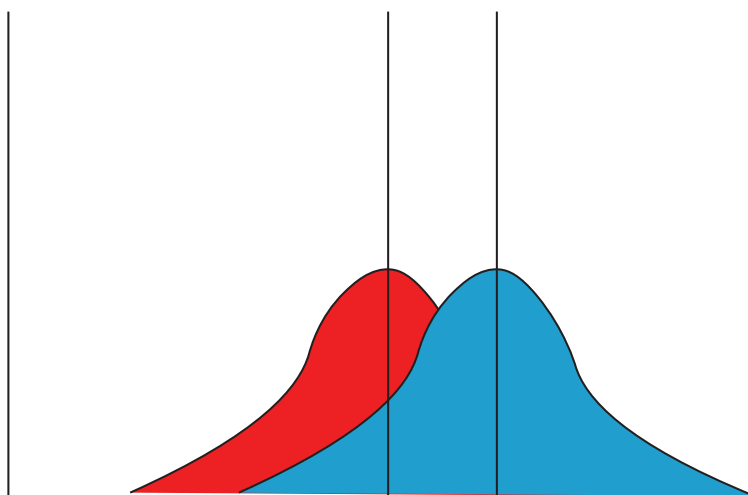
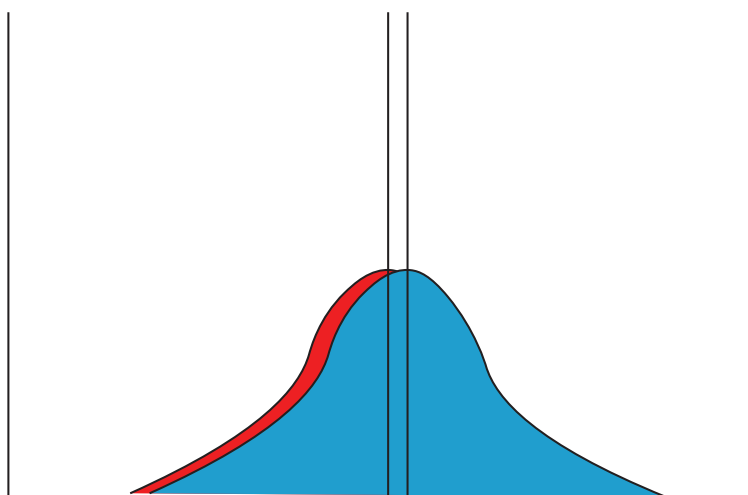
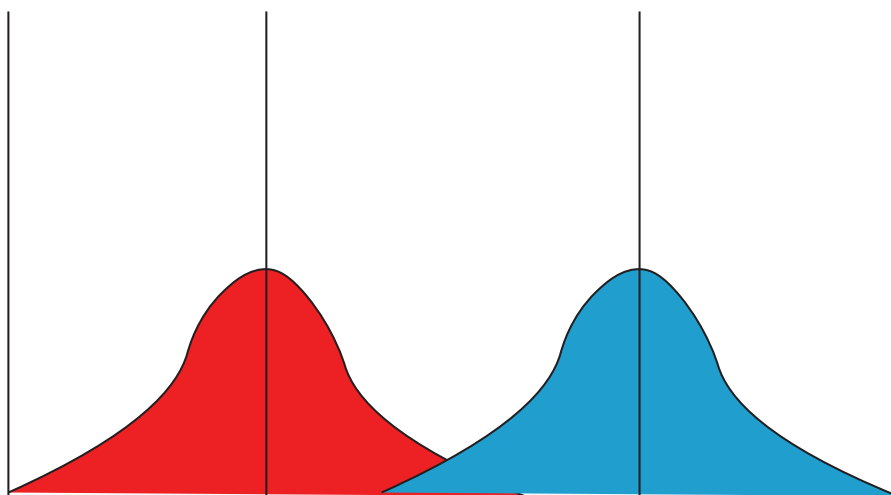
Ok, so stepping back - what is the fundamental question we're asking here?

Pose this to the class.

**Fundamentally, we want to know if there a difference between the two manufacturers in these two variables - battery life and signal distance**

They look different in the graphs you made, but is this enough? What do I mean?

Draw three different histograms on the board.



How do we tell whether there is an **actual** or **meaningful** difference in the means?

## Introduction to a t-test

There is a formalized set of statistical tools to answer this question. Some of you may have heard of t-tests before - which is what we use on this type of data.

Let's briefly talk through the logic here.

1. Our data are a sample of a larger population (think about all of the collars ever produced for both companies).
2. What we're really interested in is the difference in the means between the two groups. If they're the same, this difference is 0. If they're different, the difference is something larger or smaller than 0.
3. the code is simple, but the output is a bit trickier:

```
t.test(battery_life ~ maker, data = collars)

##
##  Welch Two Sample t-test
##
## data:  battery_life by maker
## t = -15.966, df = 89.015, p-value < 2.2e-16
## alternative hypothesis: true difference in means between group Budget Collars LLC and group Collarium Inc.
## 95 percent confidence interval:
##  -38.98179 -30.35307
## sample estimates:
## mean in group Budget Collars LLC      mean in group Collarium Inc.
##                86.79449                121.46192
```

4. Breaking down the information here:

- a. *t*: is the test statistic - here, it's a metric of how different the difference of the two means is from 0. High (or low) values indicate a big difference, small values are closer to 0.
- b. *df*: is degrees of freedom - it's a measure of your sample size and some other stuff - not something we're really going to focus on here.
- c. *p-value*: a measure of certainty of the difference outlined in the *t-statistic* above, it is veryyyyyy small. This means that there is an extremely low probability that obtaining a difference in means of 0 is very, very unlikely. There are a lot of benchmarks in different fields for what this value should be below to consider the difference *significant*, typically  $< 0.05$  is considered significant.
- d. *95 percent confidence interval*: this is the range that we can expect the test statistic (difference in means) to fall in 95% of the time given the data.

## Exploring a t-test on your own

Work on your own (but feel free to consult your group) to run a t-test on the other variable we are interested in (signal distance). Be able to describe each piece of the output, what it means, and the overall finding for the test. Also be able to relate it back to the figures we generated above.