

# ENVS 193 - HW 4

Riley Zamora

2023-05-23

## Load data using here package

```
library(here)

## here() starts at C:/Users/riley/Documents/ENVS 193/ENVS-193DS_homework-04_Zamora-Riley
data_path <- here('ntl6_v12.csv') # returns the path to specified file
data <- read.csv(data_path) # load in data specified by here function
head(data)
```

	lakeid	year4	sampldate	gearid	spname	sampletype	depth	rep	indid	length
## 1	AL	1981	1981-08-04	BSEINE	BLACKCRAPPIE		-1	4	321M1	42
## 2	AL	1981	1981-08-04	BSEINE	BLACKCRAPPIE		-1	4	321M2	41
## 3	AL	1981	1981-08-04	BSEINE	BLACKCRAPPIE		-1	4	321M3	46
## 4	AL	1981	1981-08-04	BSEINE	BLACKCRAPPIE		-1	4	321M4	26
## 5	AL	1981	1981-08-04	BSEINE	BLACKCRAPPIE		-1	4	321M5	21
## 6	AL	1981	1981-08-04	BSEINE	BLACKCRAPPIE		-1	4	321R1	240

	weight	sex	fishpart	spseq	flag
## 1	NA			3211010-100	
## 2	NA			3211010-100	
## 3	NA			3211010-100	
## 4	NA		K	3211010-100	
## 5	NA		K	3211010-100	
## 6	210		S	3211010-100	

## Problem 1

Question at hand:

**How does fish length predict fish weight for trout perch (across all sample years)?**

Our model (SLR):

Let  $Y$  represent the weight for all sampled trout perch.

Let  $x_1$  represent the length for all sampled trout perch.

Let  $\beta_0$  be the intercept coefficient.

Let  $\beta_1$  be the slope coefficient.

Let  $\epsilon$  be the error term.

$$Y = \beta_0 + \beta_1 x_1 + \epsilon$$

## Wrangle data

```
library(tidyverse)
data_wrangle <- data %>%
  select(spname, length, weight, year4) %>%
  filter(spname == 'TROUTPERCH')
```

### Part 1.

Null hypothesis: There is no linear relationship between fish length and fish weight for trout perch across all sample years.

$$H_0 : \beta_1 = 0$$

Alternate hypothesis: There is a linear relationship between fish length and fish weight for trout perch across all sample years.

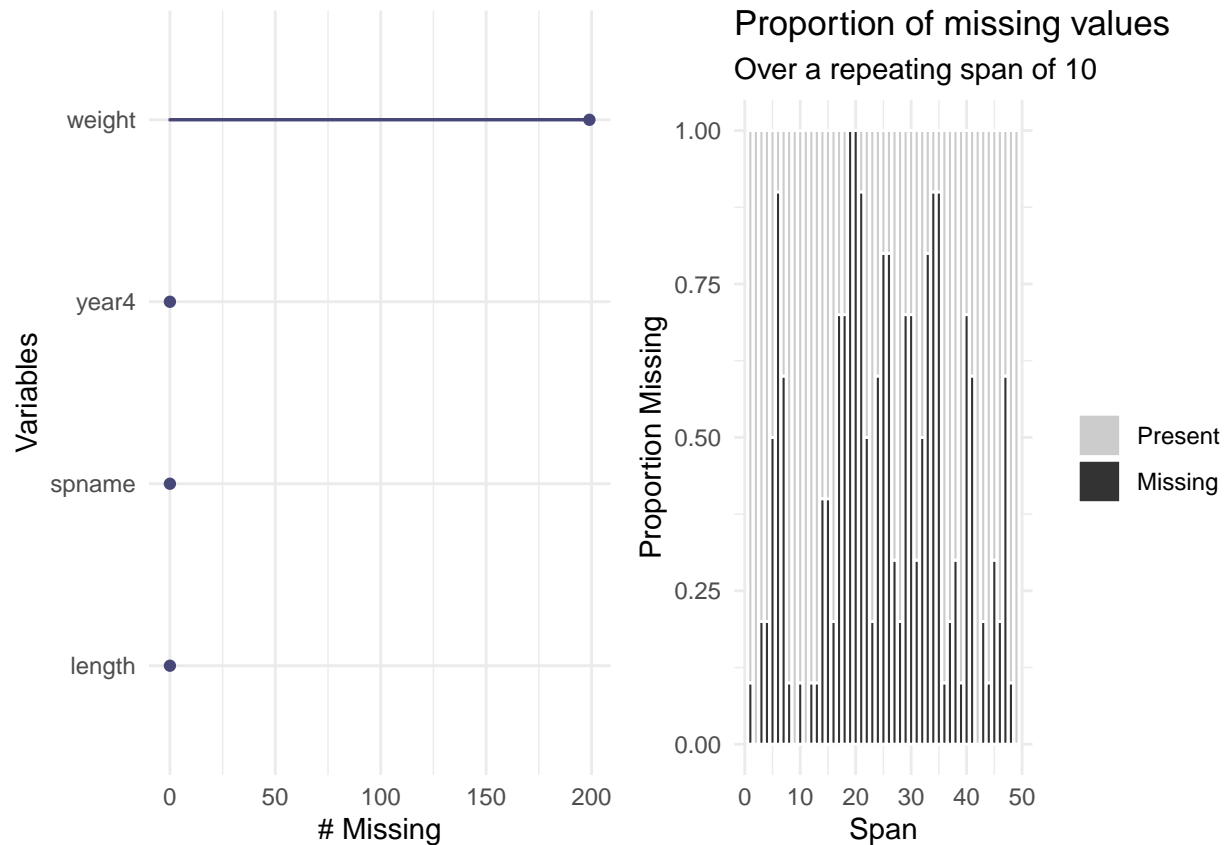
$$H_a : \beta_1 \neq 0$$

### Part 2.

```
library(naniar) # for missing data vis
library(gridExtra) # to help arrange my plots

##
## Attaching package: 'gridExtra'
## The following object is masked from 'package:dplyr':
##
##      combine
missing_var <- gg_miss_var(data_wrangle)
weight_span <- gg_miss_span(data_wrangle, weight, span_every = 10)

grid.arrange(missing_var, weight_span, layout_matrix = rbind(c(1,2)))
```



### Sub-part A:

In these missing value visualizations, we can see exactly what variables contain missing values and the span of how much missing data is in the variables that contain missing data. The left hand plot shows that **weight** is the only variable in my filtered data set. The right hand gives a better idea of the proportion of missing data for a span of 10 data points. This can help us understand if we have bias in the data and also give us an idea of whether some assumptions may be violated. We can also use these visualizations to help us decide how we can handle our missing values.

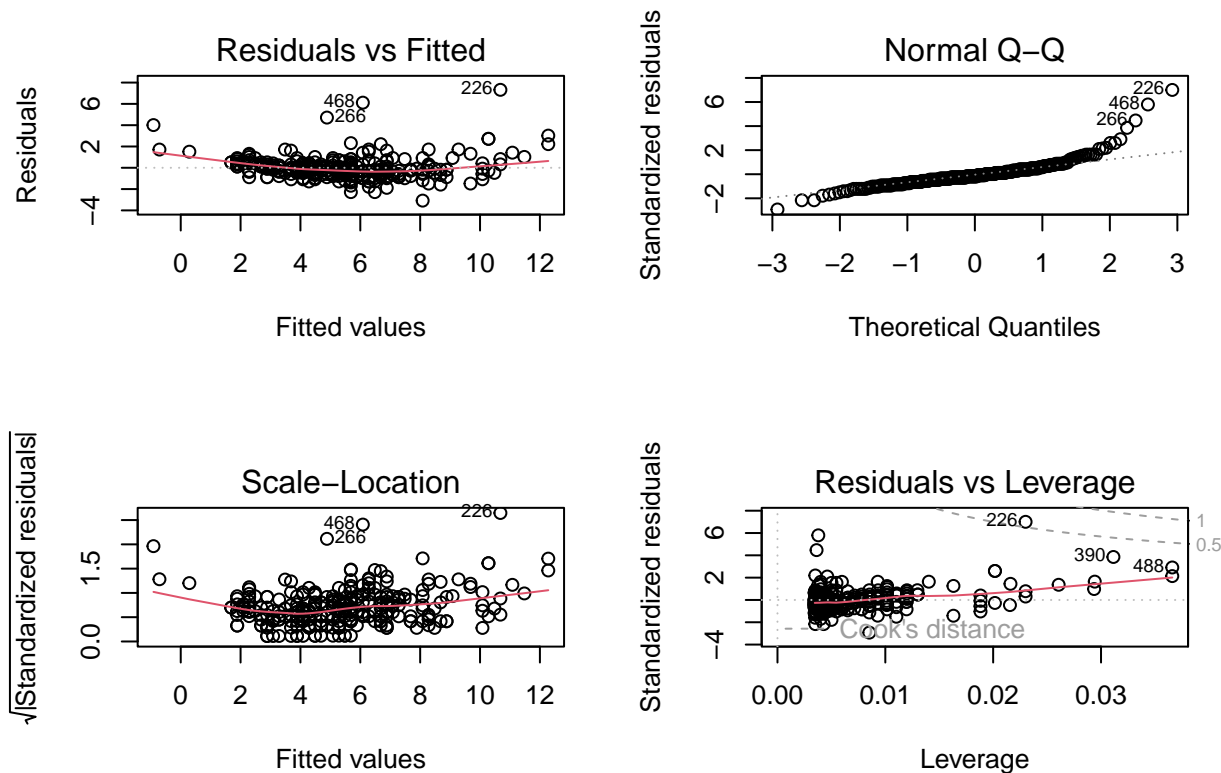
### Part 3.

```
fit <- lm(weight~length, data = data_wrangle) # fit a linear model
fit

##
## Call:
## lm(formula = weight ~ length, data = data_wrangle)
##
## Coefficients:
## (Intercept)      length
##      -11.7025       0.1999
```

### Part 4.

```
par(mfrow = c(2,2)) # aligns our resulting plots in a grid
plot(fit) # shows the assumption plots
```



## Part 5.

**Residuals Vs Fitted:** This plot is showing us the *predicted* values vs the residuals. In other words, we are seeing what values our model gives and the residuals which are the difference between the actual values and predicted value. As a rule of thumb, if the smoothed curve is relatively flat then we pass the linearity assumption. I am deciding that we have some outliers but in general we pass linearity assumption.

**Normal Q-Q:** This plot shows us standardized residuals vs theoretical quantiles of  $N(0, 1)$ . In other words it tells us if  $Y$  is normally distributed. If the points falls closely to the line, then we can say the data is normal and in our case it does.

**Scale-Location:** This plot shows  $\sqrt{\text{standardized residuals}}$  vs Fitted values which tells us if we have constant variance. Since the variability around the red line is even on either side and the red line is relatively flat for the amount of data we have, we pass the homoscedasticity assumption.

**Residuals Vs Leverage:** This plot shows us our standardized residuals vs the leverage points which helps us see if there are influential points that alter our data. We have some outliers and some and some leverage points that seem to slightly affect our data.

## Part 6.

```
summary(fit) # the summary function will show all the estimates and coefficients
```

```
##
## Call:
## lm(formula = weight ~ length, data = data_wrangle)
```

```
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.0828 -0.4862 -0.1830  0.4128  7.3191
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -11.702476   0.481564  -24.30  <2e-16 ***
## length       0.199852   0.005584   35.79  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.057 on 288 degrees of freedom
## (199 observations deleted due to missingness)
## Multiple R-squared:  0.8164, Adjusted R-squared:  0.8158
## F-statistic: 1281 on 1 and 288 DF, p-value: < 2.2e-16
```

## Part 7.