

Week 2 - Homework

STAT 420, Summer 2017, Dalpiaz

Exercise 1 (Using `lm`)

For this exercise we will use the `cats` dataset from the `MASS` package. You should use `?cats` to learn about the background of this dataset.

(a) Suppose we would like to understand the size of a cat's heart based on the body weight of a cat. Fit a simple linear model in R that accomplishes this task. Store the results in a variable called `cat_model`. Output the result of calling `summary()` on `cat_model`.

Solution:

```
library(MASS)
cat_model = lm(Hwt ~ Bwt, data = cats)
summary(cat_model)

##
## Call:
## lm(formula = Hwt ~ Bwt, data = cats)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.5694 -0.9634 -0.0921  1.0426  5.1238
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -0.3567      0.6923  -0.515   0.607
## Bwt           4.0341      0.2503  16.119 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.452 on 142 degrees of freedom
## Multiple R-squared:  0.6466, Adjusted R-squared:  0.6441
## F-statistic: 259.8 on 1 and 142 DF,  p-value: < 2.2e-16
```

(b) Output only the estimated regression coefficients. Interpret β_0 and $\hat{\beta}_1$ in the *context of the problem*. Be aware that only one of those is an estimate.

Solution:

```
coef(cat_model)

## (Intercept)      Bwt
## -0.3566624    4.0340627
```

- The mean cat heart weight (in grams) for a bodyweight of 0 kg is β_0 .
- For each additional kg of bodyweight, the estimated mean heart weight increases by $\hat{\beta}_1 = 4.0340627$ grams.

(c) Use your model to predict the heart weight of a cat that weights **3.3** kg. Do you feel confident in this prediction? Briefly explain.

Solution:

```
range(cats$Bwt)
```

```
## [1] 2.0 3.9
```

```
predict(cat_model, data.frame(Bwt = 3.3))
```

```
##          1
```

```
## 12.95574
```

Since 3.3 is inside the data range, we are interpolating, thus we are somewhat confident in our prediction. We will see later how to put bounds around this value that quantify our certainty.

(d) Use your model to predict the heart weight of a cat that weights **1.5** kg. Do you feel confident in this prediction? Briefly explain.

Solution:

```
predict(cat_model, data.frame(Bwt = 1.5))
```

```
##          1
```

```
## 5.694432
```

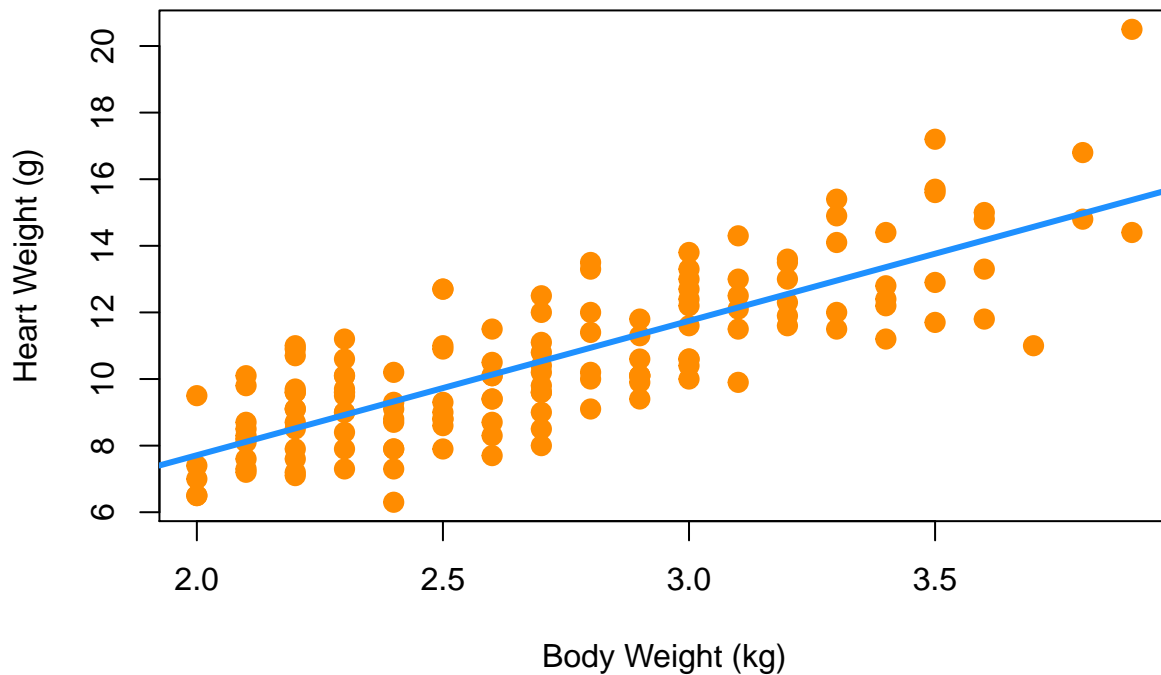
Since 1.5 is outside the data range, we are extrapolating, so we are **not** confident in our prediction. We will be able to put bounds on this, like the previous prediction; however, they will be rather meaningless.

(e) Create a scatterplot of the data and add the fitted regression line. Make sure your plot is well labeled and is somewhat visually appealing.

Solution:

```
plot(Hwt ~ Bwt, data = cats,
     xlab = "Body Weight (kg)",
     ylab = "Heart Weight (g)",
     main = "Cats: Heart Size vs Bodyweight",
     pch = 20,
     cex = 2,
     col = "darkorange")
abline(cat_model, lwd = 3, col = "dodgerblue")
```

Cats: Heart Size vs Bodyweight



(f) Report the value of R^2 for the model. Do so directly. Do not simply copy and paste the value from the full output in the console after running `summary()` in part (a).

Solution:

```
summary(cat_model)$r.squared
```

```
## [1] 0.6466209
```

This model has $R^2 = 0.6466$ which means that 64.66% of the observed variation in heart weight is explained by the linear relationship with body weight.

Exercise 2 (Writing Functions)

This exercise is a continuation of Exercise 1.

(a) Write a function called `get_sd_est` that calculates an estimate of σ in one of two ways depending on input to the function. The function should take two arguments as input:

- `model_resid` - A vector of residual values from a fitted model.
- `mle` - A logical (TRUE / FALSE) variable which defaults to FALSE.

The function should return a single value:

- s_e if `mle` is set to FALSE.
- $\hat{\sigma}$ if `mle` is set to TRUE.

Solution:

```
get_sd_est = function(model_resid, mle = FALSE) {
  n = length(model_resid) - 2 * !mle
  sqrt(sum(model_resid ^ 2) / n)
}
```

(b) Run the function `get_sd_est` on the residuals from the model in Exercise 1, with `mle` set to `FALSE`.

Solution:

```
get_sd_est(resid(cat_model))
```

```
## [1] 1.452373
```

$$s_e = 1.452373$$

(c) Run the function `get_sd_est` on the residuals from the model in Exercise 1, with `mle` set to `TRUE`.

Solution:

```
get_sd_est(resid(cat_model), mle = TRUE)
```

```
## [1] 1.442252
```

$$\hat{\sigma} = 1.442252$$

(d) To check your work, output `summary(cat_model)$sigma`. It should match at least one of (b) or (c).

Solution:

```
summary(cat_model)$sigma
```

```
## [1] 1.452373
```

We see that this value matches s_e .

Exercise 3 (Simulating SLR)

Consider the model

$$Y_i = -4 + 2x_i + \epsilon_i$$

with

$$\epsilon_i \sim N(\mu = 0, \sigma^2 = 6.25)$$

where $\beta_0 = -4$ and $\beta_1 = 2$.

This exercise relies heavily on generating random observations. To make this reproducible we will set a seed for the randomization. Alter the following code to make `birthday` store your birthday in the format: `yyyymmdd`. For example, William Gosset, better known as *Student*, was born on June 13, 1876, so he would use:

```
birthday = 18760613
set.seed(birthday)
```

(a) Use R to simulate `n = 50` observations from the above model. For the remainder of this exercise, use the following “known” values of x .

```
x = runif(n = 50, 0, 10)
```

You may use the `sim_slr` function provided in the text. Store the data frame this function returns in a variable of your choice. Note that this function calls y `response` and x `predictor`.

Solution:

```
sim_slr = function(x, beta_0 = 10, beta_1 = 5, sigma = 1) {  
  n = length(x)  
  epsilon = rnorm(n, mean = 0, sd = sigma)  
  y = beta_0 + beta_1 * x + epsilon  
  data.frame(predictor = x, response = y)  
}  
  
sim_data = sim_slr(x = x, beta_0 = -4, beta_1 = 2, sigma = 2.5)
```

(b) Fit a model to your simulated data. Report the estimated coefficients. Are they close to what you would expect? Briefly explain.

Solution:

```
sim_model = lm(response ~ predictor, data = sim_data)  
coef(sim_model)
```

```
## (Intercept)    predictor  
##   -3.355393     1.861153
```

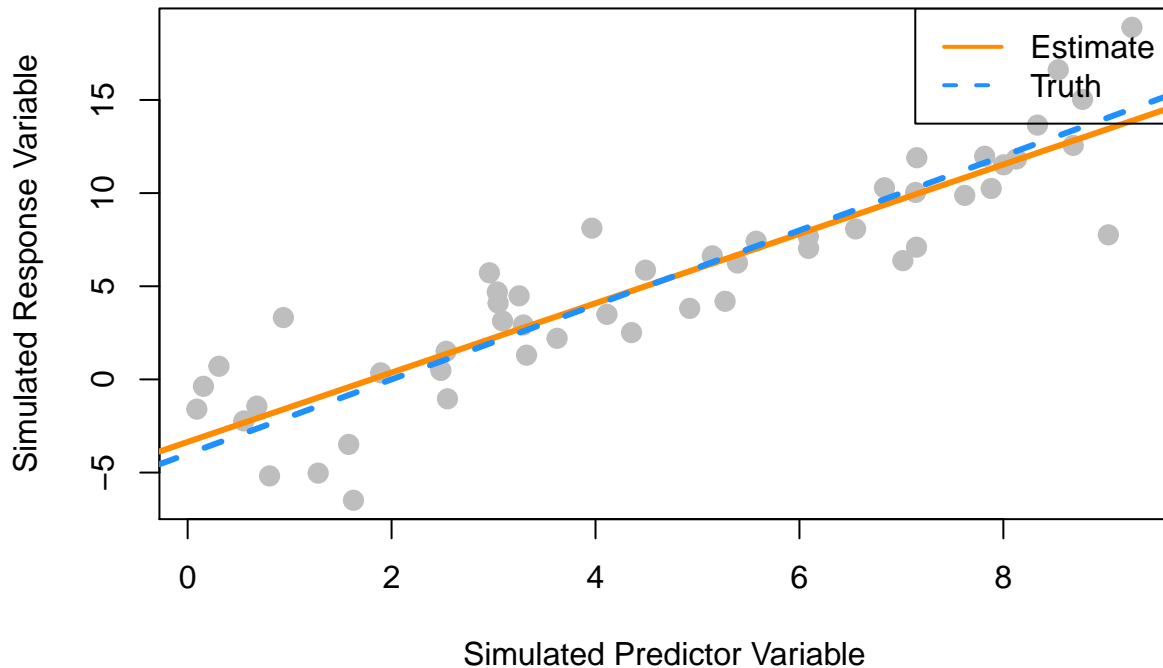
Yes, they are somewhat close to the known parameters of the model, $\beta_0 = -4$ and $\beta_1 = 2$.

(c) Plot the data you simulated in part (a). Add the regression line from part (b) as well as the line for the true model. Hint: Keep all plotting commands in the same chunk.

Solution:

```
plot(response ~ predictor, data = sim_data,  
      xlab = "Simulated Predictor Variable",  
      ylab = "Simulated Response Variable",  
      main = "Simulated Regression Data",  
      pch = 20,  
      cex = 2,  
      col = "grey")  
abline(sim_model, lwd = 3, lty = 1, col = "darkorange")  
abline(-4, 2, lwd = 3, lty = 2, col = "dodgerblue")  
legend("topright", c("Estimate", "Truth"), lty = c(1, 2), lwd = 2,  
      col = c("darkorange", "dodgerblue"))
```

Simulated Regression Data



(d) Use R to repeat the process of simulating $n = 50$ observations from the above model 2000 times. Each time fit a SLR model to the data and store the value of $\hat{\beta}_1$ in a variable called `beta_hat_1`. Some hints:

- Consider a `for` loop.
- Create `beta_hat_1` before writing the `for` loop. Make it a vector of length 2000 where each element is 0.
- Inside the body of the `for` loop, simulate new y data each time. Use a variable to temporarily store this data together with the known x data as a data frame.
- After simulating the data, use `lm()` to fit a regression. Use a variable to temporarily store this output.
- Use the `coef()` function and `[]` to extract the correct estimated coefficient.
- Use `beta_hat_1[i]` to store in elements of `beta_hat_1`.
- See the notes on Distribution of a Sample Mean for some inspiration.

You can do this differently if you like. Use of these hints is not required.

Solution:

```
beta_hat_1 = rep(0, 2000)
for(i in 1:2000) {
  sim_data = sim_slr(x = x, beta_0 = -4, beta_1 = 2, sigma = 2.5)
  model = lm(response ~ predictor, data = sim_data)
  beta_hat_1[i] = coef(model)[2]
}
```

(e) Report the mean and standard deviation of `beta_hat_1`. Do either of these look familiar?

Solution:

```
c(mean(beta_hat_1), sd(beta_hat_1))
```

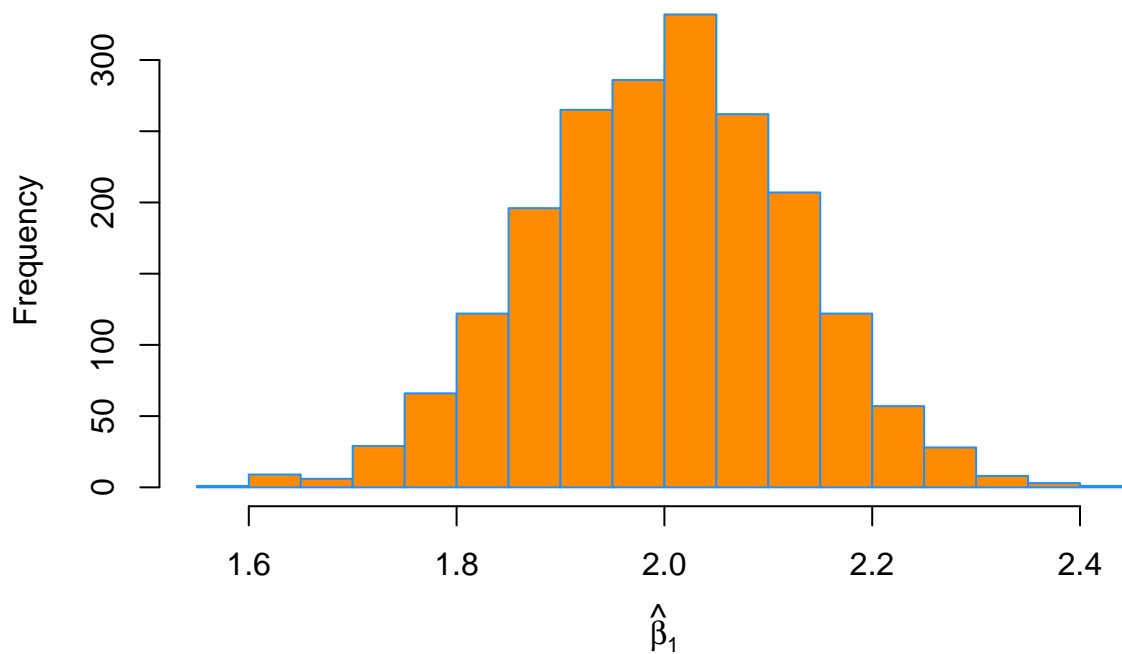
```
## [1] 1.9995458 0.1238084
```

The mean of `beta_hat_1` appears very similar to the parameter $\beta_1 = 2$.

(f) Plot a histogram of `beta_hat_1`. Comment on the shape of this histogram.

Solution:

```
hist(beta_hat_1,
     main = "",
     xlab = expression(hat(beta)[1]),
     col = "darkorange",
     border = "dodgerblue")
```



This histogram appears to be that of a normal distribution with a mean close to 2.

Exercise 4 (Be a Skeptic)

Consider the model

$$Y_i = 10 + 0x_i + \epsilon_i$$

with

$$\epsilon_i \sim N(\mu = 0, \sigma^2 = 1)$$

where $\beta_0 = 10$ and $\beta_1 = 0$.

Before answering the following parts, set a seed value equal to **your** birthday, as was done in the previous exercise.

```
birthday = 18760613
set.seed(birthday)
```

(a) Use R to repeat the process of simulating $n = 25$ observations from the above model 1500 times. For the remainder of this exercise, use the following “known” values of x .

```
x = runif(n = 25, 0, 10)
```

Each time fit a SLR model to the data and store the value of $\hat{\beta}_1$ in a variable called `beta_hat_1`. You may use the `sim_slr` function provided in the text. Hint: Yes $\beta_1 = 0$.

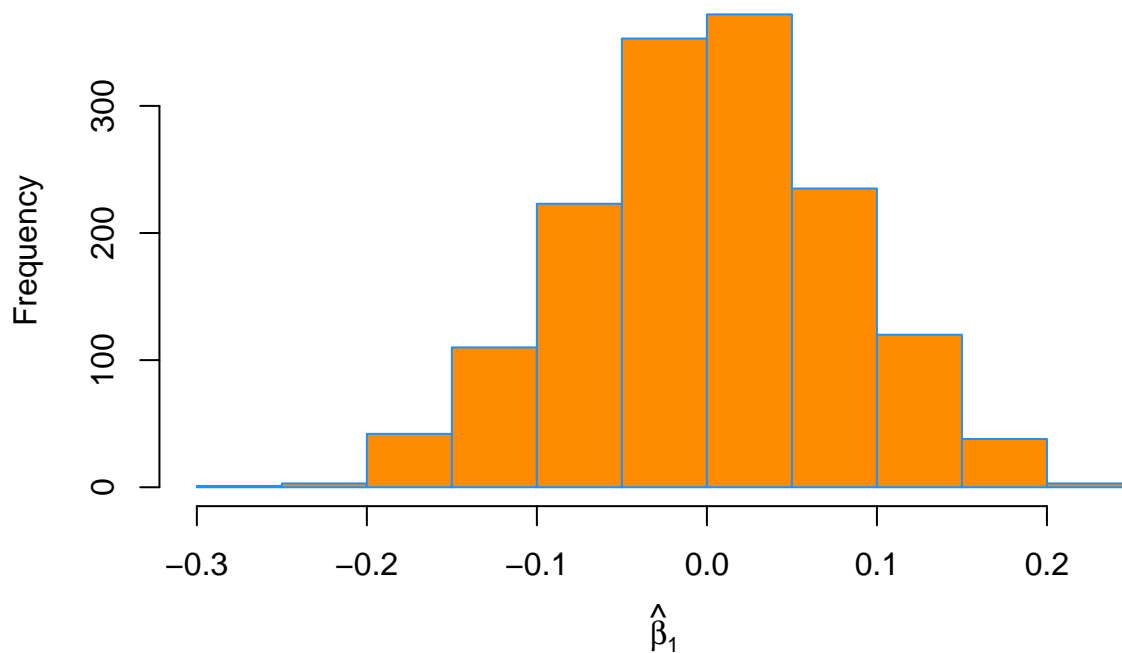
Solution:

```
beta_hat_1 = rep(0, 1500)
for(i in seq_along(beta_hat_1)) {
  sim_data = sim_slr(x = x, beta_0 = 10, beta_1 = 0, sigma = 1)
  model = lm(response ~ predictor, data = sim_data)
  beta_hat_1[i] = coef(model)[2]
}
```

(b) Plot a histogram of `beta_hat_1`. Comment on the shape of this histogram.

Solution:

```
hist(beta_hat_1,
     main = "",
     xlab = expression(hat(beta)[1]),
     col = "darkorange",
     border = "dodgerblue")
```

This histogram appears to be that of a normal distribution with a mean close to 0.

(c) Import the data in `skeptic.csv` and fit a SLR model. The variable names in `skeptic.csv` follow the same convention as those returned by `sim_slr()`. Extract the fitted coefficient for β_1 .

Solution:

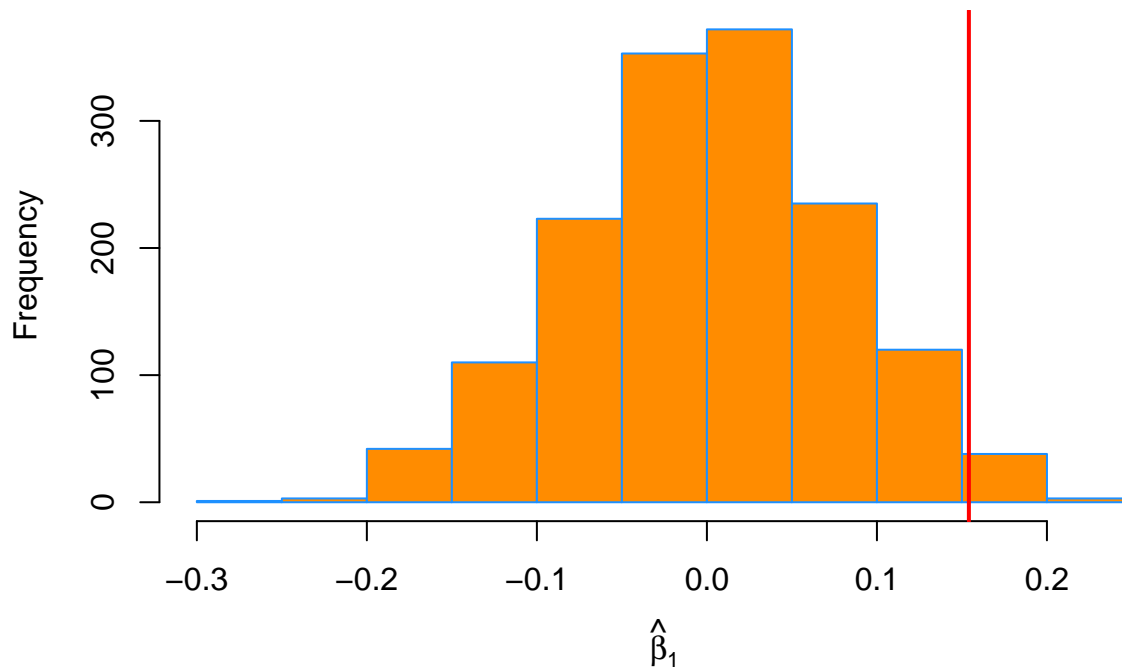
```
skeptic = read.csv("skeptic.csv")
skeptic_model = lm(response ~ predictor, data = skeptic)
coef(skeptic_model)[2]
```

```
## predictor
## 0.154081
```

(d) Re-plot the histogram from (b). Now add a vertical red line at the value of $\hat{\beta}_1$ in part (c). To do so, you'll need to use `abline(v = c, col = "red")` where `c` is your value.

Solution:

```
hist(beta_hat_1,
     main = "",
     xlab = expression(hat(beta)[1]),
     col = "darkorange",
     border = "dodgerblue")
abline(v = coef(skeptic_model)[2], col = "red", lwd = 2)
```



(e) Your value of $\hat{\beta}_1$ in (c) should be positive. What proportion of the `beta_hat_1` values are larger than your $\hat{\beta}_1$? Return this proportion, as well as this proportion multiplied by 2.

Solution:

```
mean(beta_hat_1 > coef(skeptic_model)[2])
```

```
## [1] 0.02666667
```

```
2 * mean(beta_hat_1 > coef(skeptic_model)[2])
```

```
## [1] 0.05333333
```

(f) Based on your histogram and part (e), do you think the `skeptic.csv` data could have been generated by the model given above? Briefly explain.

Solution:

```
range(beta_hat_1)
```

```
## [1] -0.2526983 0.2235652
```

While it is certainly possible, since 0.154081 is within the range of $\hat{\beta}_1$ values that we simulated, it is one of the more extreme values, with only 2.67% of the simulated values being larger. So, possible? Yes. Probable? Not exactly.

Exercise 5 (Comparing Models)

For this exercise we will use the data stored in `goalies.csv`. It contains career data for all 716 players in the history of the National Hockey League to play goaltender through the 2014-2015 season. The variables in the dataset are:

- **Player** - NHL Player Name
- **First** - First year of NHL career
- **Last** - Last year of NHL career
- **GP** - Games Played
- **GS** - Games Started
- **W** - Wins
- **L** - Losses
- **TOL** - Ties/Overtime/Shootout Losses
- **GA** - Goals Against
- **SA** - Shots Against
- **SV** - Saves
- **SV_PCT** - Save Percentage
- **GAA** - Goals Against Average
- **SO** - Shutouts
- **MIN** - Minutes
- **G** - Goals (that the player recorded, not opponents)
- **A** - Assists (that the player recorded, not opponents)
- **PTS** - Points (that the player recorded, not opponents)
- **PIM** - Penalties in Minutes

For this exercise we will define the “Root Mean Square Error” of a model as

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}.$$

(a) Fit three SLR models, each with “wins” as the response. For the predictor, use “minutes”, “goals against”, and “shutouts” respectively. For each, calculate RMSE and R^2 . Arrange the results in a markdown table, with a row for each model. Suggestion: Create a data frame that stores the results, then investigate the `kable()` function from the `knitr` package.

Solution:

```
# read data
goalies = read.csv("goalies.csv")

# fit models
goalies_model_min = lm(W ~ MIN, data = goalies)
goalies_model_ga = lm(W ~ GA, data = goalies)
goalies_model_so = lm(W ~ SO, data = goalies)

# helper functions
get_rmse = function(model) {
  sqrt(mean(resid(model) ^ 2))
}

get_r2 = function(model) {
  summary(model)$r.squared
}
```

```

# create list of model
model_list = list(goalies_model_min, goalies_model_ga, goalies_model_so)

# obtain results
results = data.frame(
  Predictor = c("`MIN`", "`GA`", "`SO`"),
  RMSE = sapply(model_list, get_rmse),
  R2 = sapply(model_list, get_r2)
)

# create markdown table of results
knitr::kable(results)

```

Predictor	RMSE	R2
MIN	16.75758	0.9711568
GA	31.08159	0.9007736
SO	44.83837	0.7934997

(b) Based on the results, which of the three predictors used is most helpful for predicting wins? Briefly explain.

Solution:

MIN since the model that uses MIN obtains the lowest RMSE, which is essentially the square root of the average of the errors of the model. So, lower is better.

Exercise 00 (SLR without Intercept)

This exercise will *not* be graded and is simply provided for your information. No credit will be given for the completion of this exercise. Give it a try now, and be sure to read the solutions later.

Sometimes it can be reasonable to assume that β_0 should be 0. That is, the line should pass through the point (0,0). For example, if a car is traveling 0 miles per hour, its stopping distance should be 0! (Unlike what we saw in the book.)

We can simply define a model without an intercept,

$$Y_i = \beta x_i + \epsilon_i.$$

(a) In the **Least Squares Approach** section of the text you saw the calculus behind the derivation of the regression estimates, and then we performed the calculation for the `cars` dataset using R. Here you need to do, but not show, the derivation for the slope only model. You should then use that derivation of $\hat{\beta}$ to write a function that performs the calculation for the estimate you derived.

In summary, use the method of least squares to derive an estimate for β using data points (x_i, y_i) for $i = 1, 2, \dots, n$. Simply put, find the value of β to minimize the function

$$f(\beta) = \sum_{i=1}^n (y_i - \beta x_i)^2.$$

Then, write a function `get_beta_no_int` that takes input:

- x - A predictor variable
- y - A response variable

The function should then output the $\hat{\beta}$ you derived for a given set of data.

Solution:

```
get_beta_no_int = function(x, y) {
  sum(x * y) / sum(x ^ 2)
}
```

(b) Write your derivation in your .Rmd file using TeX. Or write your derivation by hand, scan or photograph your work, and insert it into the .Rmd as an image. See the RMarkdown documentation for working with images.

Solution:

To minimize the function

$$f(\beta) = \sum_{i=1}^n (y_i - \beta x_i)^2$$

we take the first derivative with respect to β

$$\frac{\partial f}{\partial \beta} = \sum_{i=1}^n 2(y_i - \beta x_i)(-x_i) = 2\beta \sum_{i=1}^n x_i^2 - 2 \sum_{i=1}^n x_i y_i.$$

Then solving

$$\frac{\partial f}{\partial \beta} = 0$$

for β gives

$$\hat{\beta} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}.$$

Lastly, we verify that

$$\frac{\partial^2 f}{\partial \beta^2} = 2 \sum_{i=1}^n x_i^2 > 0.$$

Thus $\hat{\beta}$ does minimize $f(\beta)$ therefor it is the least-squares estimator of β .

(c) Test your function on the `cats` data using body weight as x and heart weight as y . What is the estimate for β for this data?

Solution:

```
get_beta_no_int(x = cats$Bwt, y = cats$Hwt)
```

```
## [1] 3.907113
```

$$\hat{\beta} = 3.9071133$$

(d) Check your work in R. The following syntax can be used to fit a model without an intercept:

```
lm(response ~ 0 + predictor, data = dataset)
```

Use this to fit a model to the `cat` data without an intercept. Output the coefficient of the fitted model. It should match your answer to (c).

Solution:

```
lm(Hwt ~ 0 + Bwt, data = cats)
```

```
##  
## Call:  
## lm(formula = Hwt ~ 0 + Bwt, data = cats)  
##  
## Coefficients:  
##      Bwt  
## 3.907
```