

Week 2 - Homework

STAT 420, Fall 2024, Banghao Chi

09/10/2024

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.

```
library(MASS)
?MASS::cats
dataset = MASS::cats
```

(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:

```
cat_model = lm(Hwt ~ Bwt, data = dataset)
summary(cat_model)
```

```
##
## Call:
## lm(formula = Hwt ~ Bwt, data = dataset)
##
## 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 $\hat{\beta}_0$ and β_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 interpretation is as follow:

- $\hat{\beta}_0$: This is the **estimated** *mean* heart weight of a cat whose body weight is zero kg.
- β_1 : This is the mean increase in *mean* heart weight of a cat for an increase in its body weight of one kg.

(c) Use your model to predict the heart weight of a cat that weights **3.1** kg. Do you feel confident in this prediction? Briefly explain.

Solution:

```
range(dataset$Bwt)
```

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

```
predicted_bwt = predict(cat_model, newdata = data.frame(Bwt = 3.1))
predicted_bwt
```

```
##          1
## 12.14893
```

The predicted heart weight of a cat that weights 3.1 kg is **12.1489319**. Since the body weight of 3.1 is within the range of observations in the dataset of 2.0 to 3.9, I feel confident in this prediction.

(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:

```
range(dataset$Bwt)
```

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

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

```
##          1
## 5.694432
```

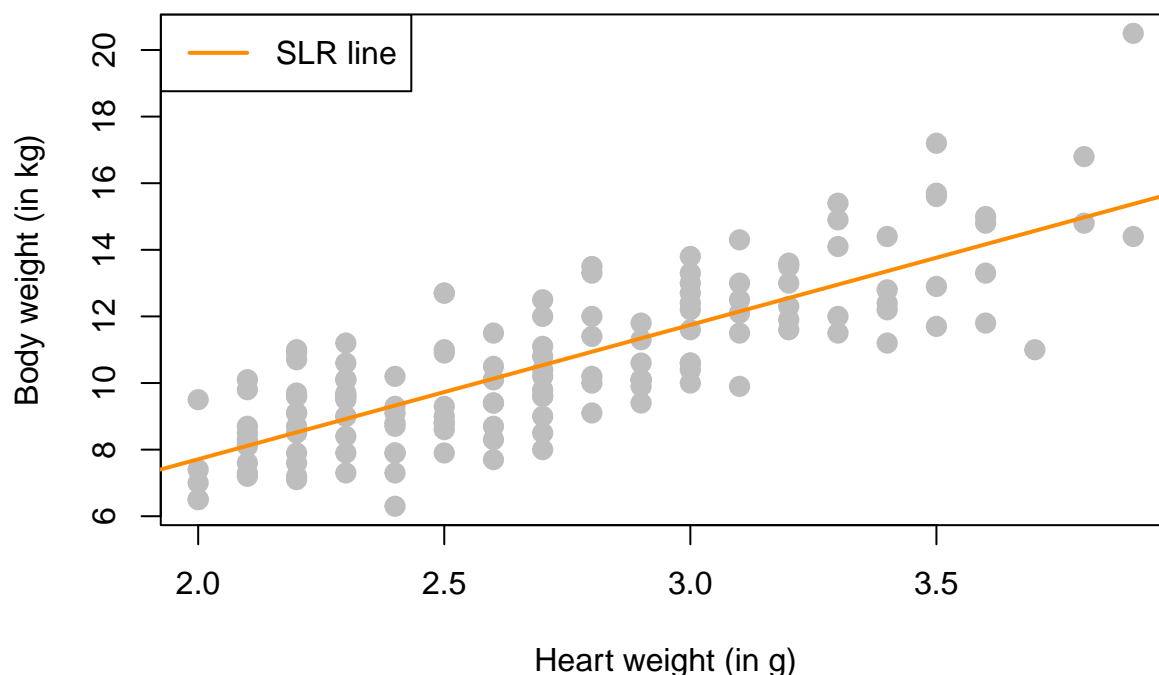
The predicted heart weight of a cat that weights 1.5 kg is **5.6944316**. Since the body weight of 1.5 is **NOT** within the range of observations in the dataset of 2.0 to 3.9, I **don't** feel confident in this prediction.

(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 = dataset,
     xlab = "Heart weight (in g)",
     ylab = "Body weight (in kg)",
     main = "Scatter plot of Heart weight v.s. Body weight of cats",
     pch = 20,
     cex = 2,
     col = "grey"
)
abline(cat_model, col = "darkorange", lwd = 2)
legend("topleft", "SLR line", col = "darkorange", lwd = 2)
```

Scatter plot of Heart weight v.s. Body weight of cats



(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:

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

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

The value of R^2 for the model is **0.6466209**.

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 three arguments as input:

- `fitted_vals` - A vector of fitted values from a model
- `actual_vals` - A vector of the true values of the response
- `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(fitted_vals, actual_vals, mle) {
  res = sum((actual_vals - fitted_vals) ^ 2)
  n = length(fitted_vals)
```

```

if (mle) {
  sqrt(res / n)
} else {
  sqrt(res / (n - 2))
}
}

```

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

Solution:

```

se_est = get_sd_est(summary(cat_model)$residuals, 0, FALSE)
se_est

```

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

This is the residual standard error, which in this context of the model means that the most of the babys' heart weights will be around a mean value with a fluctuation of ± 1.452373 or even ± 2.9047467 , which in this case is an unbiased estimation of the sigma.

(c) Run the function `get_sd_est` on the residuals from the model in Exercise 1, with `mle` set to `TRUE`. Explain the resulting estimate in the context of the model. Note that we are trying to estimate the same parameter as in part (b).

Solution:

```

sd_est = get_sd_est(summary(cat_model)$residuals, 0, TRUE)
sd_est

```

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

This is the residual standard error computed based on maximum likelihood estimate, which in this context of the model means that the most of the babys' heart weights will be around a mean value with a fluctuation of ± 1.442252 or even ± 2.8845043 grams, which in this case is a biased estimation of the sigma.

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

Solution:

```

resid_summary = summary(cat_model)$residuals
sigma = summary(cat_model)$sigma
sigma

```

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

```
sigma == get_sd_est(resid_summary, 0, FALSE) | sigma == get_sd_est(resid_summary, 0, TRUE)
```

```
## [1] TRUE
```

Exercise 3 (Simulating SLR)

Consider the model

$$Y_i = 5 + -3x_i + \epsilon_i$$

with

$$\epsilon_i \sim N(\mu = 0, \sigma^2 = 10.24)$$

where $\beta_0 = 5$ and $\beta_1 = -3$.

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 = 20021023
set.seed(birthday)
```

(a) Use R to simulate `n = 25` observations from the above model. For the remainder of this exercise, use the following “known” values of x .

```
x = runif(n = 25, 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, beta_0 = 5, beta_1 = -3, sigma = sqrt(10.24))
head(sim_data, 5)
```

```
## predictor response
## 1 9.6335418 -28.194498
## 2 9.1504591 -20.605099
## 3 0.5035793  2.139580
## 4 6.8909761 -14.876118
## 5 2.6483257  -3.919906
```

(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
##      3.329420    -2.756679
```

The coefficients are close to what I would expect, since **3.329420** is close to 5 and **-2.756679** is close to -3.

(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.

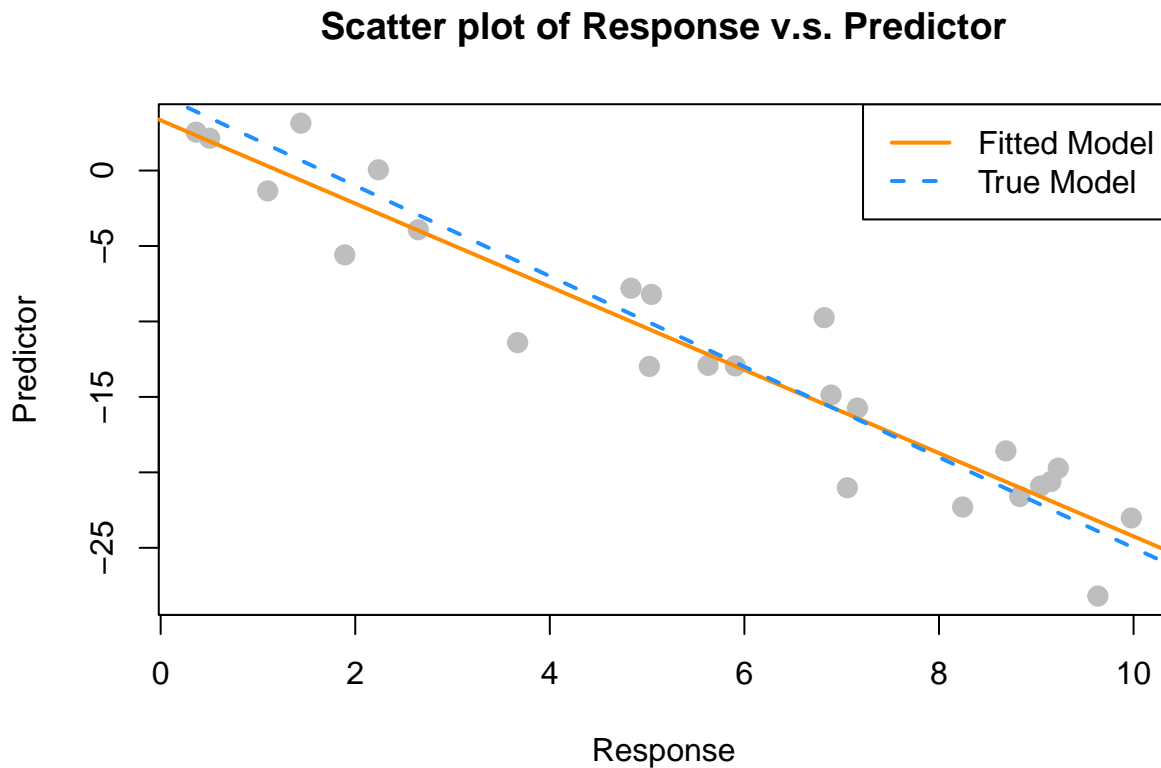
Solution:

```
plot(response ~ predictor, data = sim_data,
     xlab = "Response",
     ylab = "Predictor",
     main = "Scatter plot of Response v.s. Predictor",
```

```

pch = 20,
cex = 2,
col = "grey"
)
abline(sim_model, col = "darkorange", lwd = 2)
abline(5, -3, col = "dodgerblue", lwd = 2, lty = 2)
legend("topright",
      legend = c("Fitted Model", "True Model"),
      col = c("darkorange", "dodgerblue"),
      lwd = 2,
      lty = c(1, 2))

```



(d) Use R to repeat the process of simulating $n = 25$ observations from the above model 1500 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 1500 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, 1500)
for (i in 1:1500) {

```

```

x = runif(n = 25, 0, 10)
sim_data = sim_slr(x, beta_0 = 5, beta_1 = -3, sigma = sqrt(10.24))
sim_model = lm(response ~ predictor, data = sim_data)
beta_hat_1[i] = coef(sim_model)[2]
}

```

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

Solution:

```

mean_beta_hat_1 = mean(beta_hat_1)
mean_beta_hat_1

```

```
## [1] -3.002742
```

```

sd_beta_hat_1 = sd(beta_hat_1)
sd_beta_hat_1

```

```
## [1] 0.2273487
```

The mean and standard deviation of `beta_hat_1` is **-3.0027419** and **0.2273487** respectively. The mean is really similar to the actual value of `beta_1` of the real model.

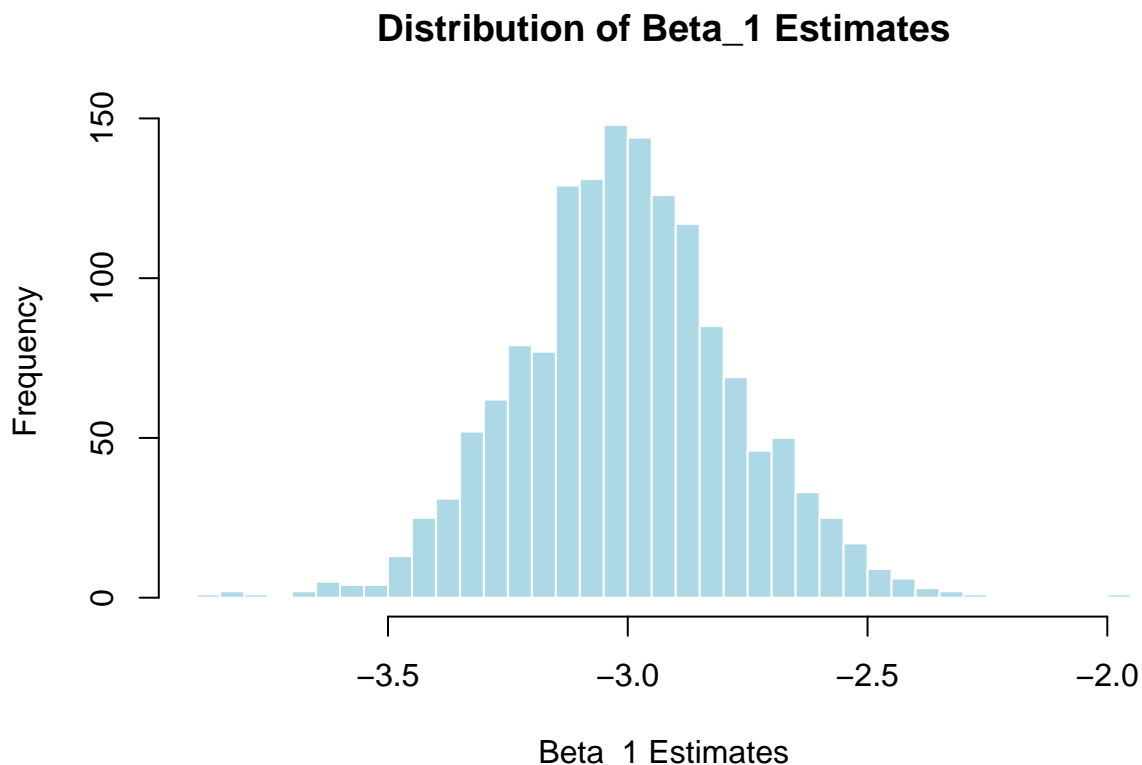
(f) Plot a histogram of `beta_hat_1`. Comment on the shape of this histogram.

Solution:

```

hist(beta_hat_1,
     col = "lightblue",
     breaks = 30,
     border = "white",
     main = "Distribution of Beta_1 Estimates",
     xlab = "Beta_1 Estimates")

```



The shape of the histogram **resembles the shape of normal distribution**, which is **symmetric**, focused on the center around the value of **-3**, and fade to 0 as the distance gets bigger.

Exercise 4 (Be a Skeptic)

Consider the model

$$Y_i = 3 + 0 \cdot x_i + \epsilon_i$$

with

$$\epsilon_i \sim N(\mu = 0, \sigma^2 = 4)$$

where $\beta_0 = 3$ 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 = 20021023
set.seed(birthday)
```

(a) Use R to repeat the process of simulating $n = 75$ observations from the above model 2500 times. For the remainder of this exercise, use the following “known” values of x .

```
x = runif(n = 75, 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:

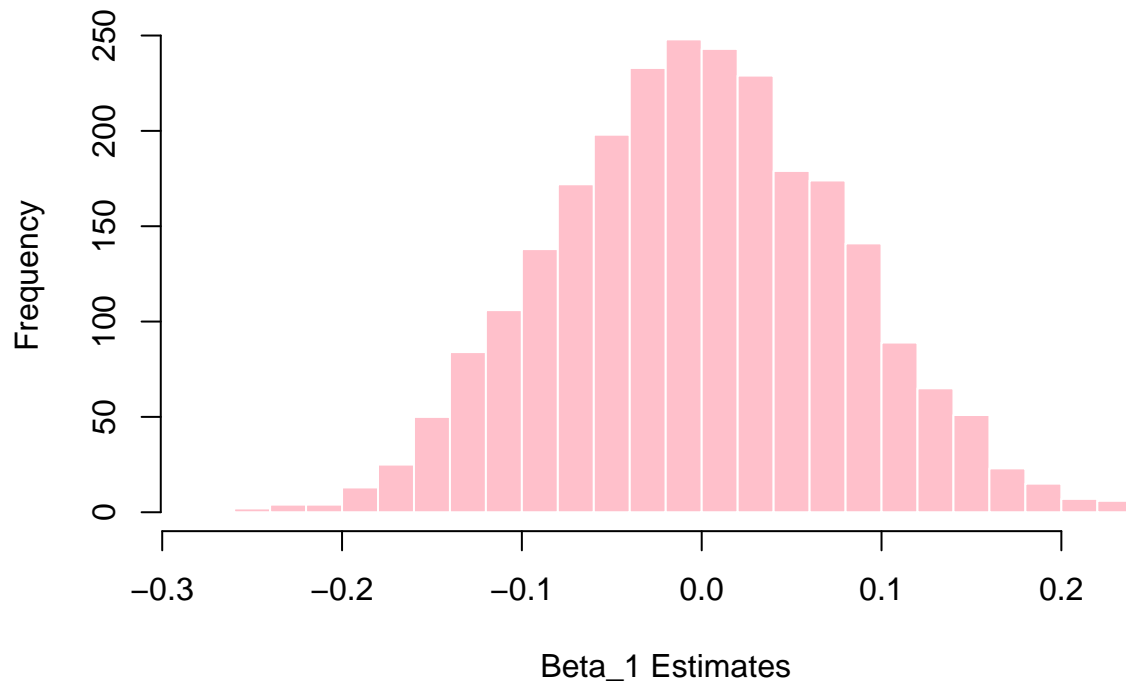
```
sim_slr = function(x, beta_0 = 10, beta_1 = 5, sigma = 1) {
  n = length(x)
  epsilon = rnorm(n, mean = 0, sd = sigma)
  y = beta_0 + beta_1 * x + epsilon
  data.frame(predictor = x, response = y)
}
beta_hat_1 = rep(0, 2500)
for (i in 1:2500) {
  x = runif(n = 75, 0, 10)
  sim_data = sim_slr(x, beta_0 = 3, beta_1 = 0, sigma = 2)
  sim_model = lm(response ~ predictor, data = sim_data)
  beta_hat_1[i] = coef(sim_model)[2]
}
```

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

Solution:

```
hist(beta_hat_1,
     col = "pink",
     breaks = 30,
     border = "white",
     main = "Distribution of Beta_1 Estimates",
     xlab = "Beta_1 Estimates")
```


Distribution of Beta_1 Estimates



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

Solution:

```
dataset = read.csv("skeptical.csv")
skeptical_model = lm(response ~ predictor, data = dataset)
beta_1_skeptical = coef(skeptical_model)[2]
beta_1_skeptical
```

```
## predictor
## -0.2221927
```

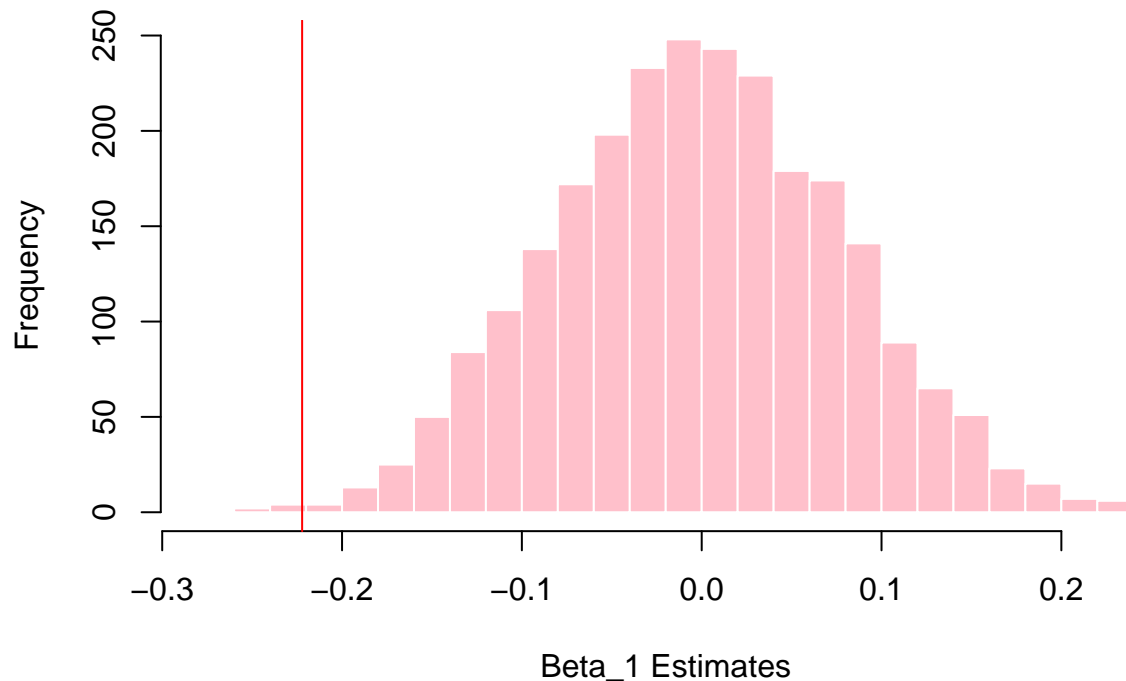
The fitted coefficient for β_1 is **-0.2221927**.

(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,
     col = "pink",
     breaks = 30,
     border = "white",
     main = "Distribution of Beta_1 Estimates",
     xlab = "Beta_1 Estimates")
abline(v = beta_1_skeptical, col = "red")
```

Distribution of Beta_1 Estimates



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

Solution:

```
proportion = mean(beta_hat_1 < beta_1_skeptic)
proportion
```

```
## [1] 0.0028
```

```
proportion_times_2 = proportion * 2
proportion_times_2
```

```
## [1] 0.0056
```

(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:

No, since:

- The proportion of `beta_hat_1` that are smaller than $\hat{\beta}_1$ is of only 0.0028, which is too small - if the data is really generated by the model given, this value should be around 0.5 since the value of $\hat{\beta}_1$ will be around the mean of `beta_hat_1`.
- From the histogram we can also draw this conclusion, since there're few frequencies that elements in `beta_hat_1` are smaller than $\hat{\beta}_1$.

Exercise 5 (Comparing Models)

For this exercise we will use the `Ozone` dataset from the `mlbench` package. You should use `?Ozone` to learn about the background of this dataset. You may need to install the `mlbench` package. If you do so, do not

include code to install the package in your R Markdown document.

For simplicity, we will perform some data cleaning before proceeding.

```
data(Ozone, package = "mlbench")
Ozone = Ozone[, c(4, 6, 7, 8)]
colnames(Ozone) = c("ozone", "wind", "humidity", "temp")
Ozone = Ozone[complete.cases(Ozone), ]
```

We have:

- Loaded the data from the package
- Subset the data to relevant variables
 - This is not really necessary (or perhaps a good idea) but it makes the next step easier
- Given variables useful names
- Removed any observation with missing values
 - This should be given much more thought in practice

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 “ozone” as the response. For the predictor, use “wind speed,” “humidity percentage,” and “temperature” 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:

```
library(knitr)

rmse = function(actual, fitted) {
  sqrt(mean((actual - fitted)^2))
}

wind_model = lm(ozone ~ wind, data = Ozone)
wind_rmse = rmse(Ozone$ozone, predict(wind_model))
wind_r_squared = summary(wind_model)$r.squared

humidity_model = lm(ozone ~ humidity, data = Ozone)
humidity_rmse = rmse(Ozone$ozone, predict(humidity_model))
humidity_r_squared = summary(humidity_model)$r.squared

temp_model = lm(ozone ~ temp, data = Ozone)
temp_rmse = rmse(Ozone$ozone, predict(temp_model))
temp_r_squared = summary(temp_model)$r.squared

results = data.frame(
  Predictor = c("Wind", "Humidity", "Temperature"),
  RMSE = c(wind_rmse, humidity_rmse, temp_rmse),
  R_squared = c(wind_r_squared,
                humidity_r_squared,
                temp_r_squared)
)
kable(results, format = "markdown", col.names = c("Predictor", "RMSE", "R-squared"))
```

Predictor	RMSE	R-squared
Wind	7.961695	0.0001402
Humidity	7.147822	0.1941105
Temperature	5.009257	0.6042011

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

Solution:

Based on the results, **temperature** is most helpful for predicting ozone readings, since:

- Its R^2 is the biggest of around **60.42%**, explaining much more variation of the data as compared to the other two factors.
- Its RMSE is the smallest compared to the other two, meaning that it has the lowest prediction error, which is most helpful for predicting ozone readings.

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.

(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.

(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?

(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).