

Lab 2

Shuyuan Wang, sw3449

September 26th, 2019

Instructions

Before the lab is due, make sure that you upload a knitted HTML or pdf file to the canvas page (this should have a .html or .pdf extension). No need to upload the .Rmd file.

Part (A): Simple Linear Regression Model

- 1) Import the **diamonds_small.csv** dataset into R and store in a dataframe called **diamonds**. Use the **lm()** command to regress **price** (response) on **carat** (predictor) and save this result as **lm0**. What are the coefficients of **lm0**? (Some of this problem is solved for you below.)

```
diamonds <- read.csv("diamonds_small.csv", as.is = TRUE, header = TRUE)
lm0 <- lm(price~carat, data=diamonds)
lm0$coefficients
```

```
## (Intercept)      carat
##   -2256.361    7756.426
```

Recall from lecture that the estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ that you just calculated with **lm()** are functions of the data values and are therefore themselves random (they inherit variability from the data). If we were to recollect the diamonds data over and over again, the estimates would be different each time.

In this lab we'll use bootstrapping to answer the following questions:

1. "How much does $\hat{\beta}_1$ vary from one replication of the experiment to the other?"
2. "What are all the values of β_1 that would have produced this data with high probability?"

Part (B): How Does $\hat{\beta}_1$ Vary?

Strategy: we'll re-sample (**price**, **carat**) pairs in order to provide an estimate for how $\hat{\beta}_1$ varies across samples.

- 1) How many rows are in the **diamonds** dataset? Call this value **n**.

```
n <- dim(diamonds)[1]
n
```

```
## [1] 53940
```

- 2) We'll next use the **sample()** function to re-sample **n** rows of the **diamonds** dataset *with replacement*. The following code provides a single re-sample of the values $1, 2, \dots, n$, or a single re-sample of the rows of the dataset.

```
resample1 <- sample(1:n, n, replace = TRUE)
```

Now write a loop to calculate **B <- 1000** such re-samples and store them as rows of the matrix **resampled_values** which will have **B** rows and **n** columns.

```
B <- 1000
resampled_values <- matrix(NA, nrow = B, ncol = n)
```

```
for (b in 1:B) {
  resampled_values[b,] <- sample(1:n, n, replace=T)
}
```

- 3) Now we'll use each re-sampled dataset to provide a new estimate of $\hat{\beta}_1$. Write a line of code that uses **resample1** above to produce a resamples dataset of (**price**, **carat**) pairs. Using the re-sampled dataset, use **lm()** to produce new estimates of $\hat{\beta}_0$ and $\hat{\beta}_1$. These values should be stored in a vector called **resample1_est**s.

```
resample1_data <- diamonds[resample1,]
resample1_est <- lm(price~carat, data=resample1_data)$coefficients
resample1_est
```

```
## (Intercept)      carat
##   -2271.225    7782.354
```

- 4) Repeat the above call for each re-sampled dataset produced from the **resampled_values** matrix. We'll store the new coefficient estimates in a matrix **resampled_est**s with **B** rows and **2** columns. Again you'll want to write a loop, this time that iterates over the rows of **resampled_values**. (Note that if you are very clever this could be done using **apply()**.) Make sure to print **head(resample_est**s) at the end.

```
resampled_est <- matrix(NA, nrow = B, ncol = 2)
colnames(resampled_est) <- c("Intercept_Est", "Slope_Est")
for (b in 1:B) {
  resampled_data <- diamonds[resampled_values[b,],]
  resampled_est[b,] <- lm(price~carat, data=resampled_data)$coefficients
}
head(resampled_est)
```

```
##      Intercept_Est Slope_Est
## [1,]    -2236.520  7723.715
## [2,]    -2287.450  7807.264
## [3,]    -2217.486  7696.825
## [4,]    -2226.212  7716.444
## [5,]    -2256.724  7747.163
## [6,]    -2270.751  7784.206
```

Hint: You may want to use multiple lines of code within the for loop. One idea is to first use the rows corresponding to re-sample **b** provided in **resampled_values** to create a resampled dataset. Then use the new dataset to provide new estimates of the coefficients.

- 5) Recall from lecture that $(\hat{\beta}_1^{(b)})_{b=1}^B - \hat{\beta}_1$ approximates the sampling distribution of $\hat{\beta}_1 - \beta_1$ where β_1 is the population parameter, $\hat{\beta}_1$ is the estimate from our original dataset, and $(\hat{\beta}_1^{(b)})_{b=1}^B$ are the B bootstrap estimates.

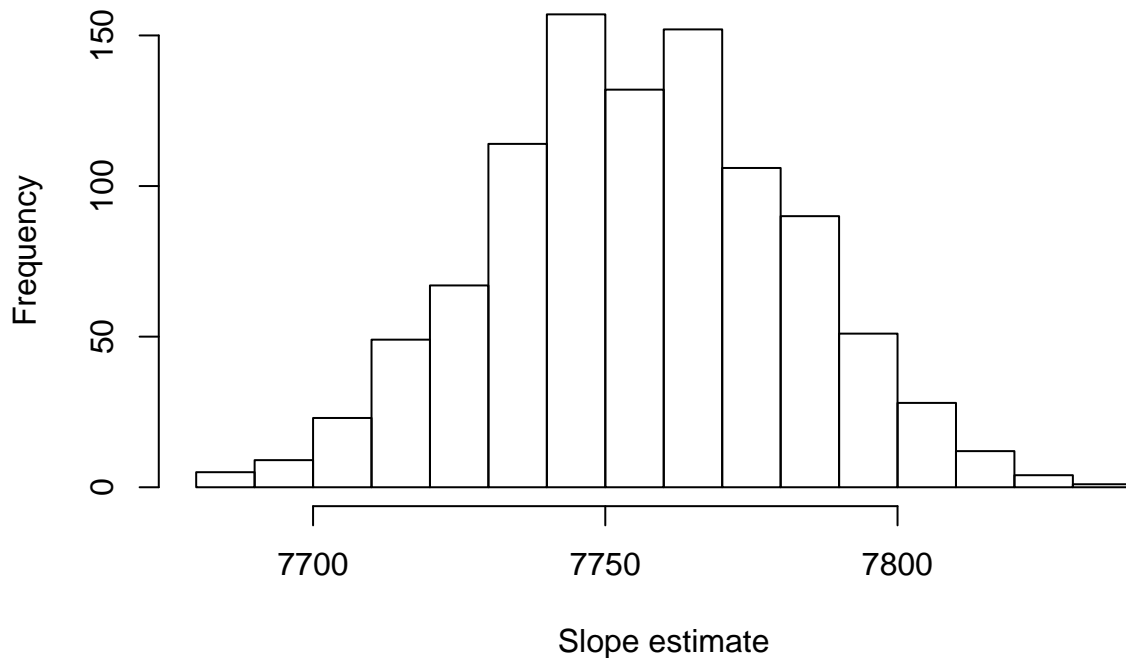
Make a vector **diff_estimates** that holds the differences between the original estimate of $\hat{\beta}_1$ from **lm0** and the bootstrap estimates. It should have length **B**.

```
diff_estimates <- resampled_est[,2] - lm0$coefficients[2]
```

- 6) Plot a histogram of the bootstrap estimates of $\hat{\beta}_1$ (they're in the 'Slope_Est' column). Label the x-axis appropriately.

```
hist(resampled_est[, 'Slope_Est'], xlab='Slope estimate', main='Slope estimates of bootstrap samples')
```

Slope estimates of bootstrap samples



7) Calculate the standard deviation of the bootstrap estimates.

```
sd(resampled_estimates[,1])
```

```
## [1] 16.45731
```

```
sd(resampled_estimates[,2])
```

```
## [1] 25.80064
```

Part (C): Bootstrap Confidence Intervals

Note: This section is optional. If you get the chance to do it during lab, great, but it's not necessary that this part is completed when you turn in the lab.

Finally we'd like to approximate confidence intervals for the regression coefficients. Recall that a confidence interval is a random interval which contains the truth with high probability (the confidence level). If the confidence interval for β_1 is C , and the confidence level is $1 - \alpha$, then we want

$$Pr(\beta_1 \in C) = 1 - \alpha$$

no matter what the true value of β_1 .

We estimate the confidence interval from the bootstrap estimates by finding a range of $(\hat{\beta}_1^{(b)})_{b=1}^B - \hat{\beta}_1$ which holds $1 - \alpha$ percent of the values. In our case, let $\alpha = 0.05$, so we estimate a confidence interval with level 0.95.

- (1) Let **Cu** and **Cl** be the upper and lower limits of the confidence interval. Use the **quantile()** function to find the 0.025 and 0.975 quantiles of the vector **diff_estimates** calculated in B(5). Then **Cu** is the sum of the original estimate of $\hat{\beta}_1$ from **lm0** with the upper quantile and **Cl** is the sum of the original estimate of $\hat{\beta}_1$ from **lm0** with the lower quantile.

```

Cl <- quantile(diff_estimates,prob=0.025)
Cu <- quantile(diff_estimates,prob=0.975)
int <- c(Cl, Cu)
int

```

```

##      2.5%      97.5%
## -50.57345  49.48705

```

- (2) Instead of traditional bootstrap intervals, construct **percentile** based bootstrap intervals. Use the **quantile()** function to find the 0.025 and 0.975 quantiles of the vector **resampled_ests[, "Slope_Est"]** calculated in B(4).

```

Cl2 <- quantile(resampled_ests[, 'Slope_Est'],prob=0.025)
Cu2 <- quantile(resampled_ests[, 'Slope_Est'],prob=0.975)
c(Cl2,Cu2)

```

```

##      2.5%      97.5%
## 7705.852 7805.913

```