

## Assignment 2

NAME: Nicholas Tolley

DUE DATE: February 28th, 6pm

### Problem 1 (100 pts)

In the earnings dataset you can find salary (*earn*) and some socio-demographic characteristics of each subject, including variables such as *height*, *weight*, gender (*male*), *ethnicity*, *education*, mother's (*mother\_education*) and father's education (*father\_education*), *walk* (e.g. walking time), *exercise*, if they smoke or not (*smokenow*), *tense*, *angry* and *age*.

The dataset can be found in Canvas in the Data folder (file name: earnings.csv):

- (a) (10 points) Subset the data and consider only the variables: *education*, *mother\_education*, *father\_education*, *walk*, *exercise*, *tense*, *angry*, *weight*, *height*. Check the correlation by performing a figure similar to Figure 1 below (make sure not to use the default colours but rather choose your own). Take special care to the labels and legend. What can you say about the results? What would you expect from a linear regression model (hint: there are some variables to be excluded/included in the model)? Perform a test statistic for the correlation between *earn* and *education*, write the hypothesis test and the results you will obtain.

```
library(ggplot2)
library(boot)
library(tidyverse)

## -- Attaching packages ----- tidyverse 1.3.2 --
## v tibble 3.1.8      v dplyr 1.1.0
## v tidyr 1.3.0      v stringr 1.5.0
## v readr 2.1.3      v forcats 1.0.0
## v purrr 1.0.1
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag() masks stats::lag()

library(ggcorrplot)
library(tidyr)

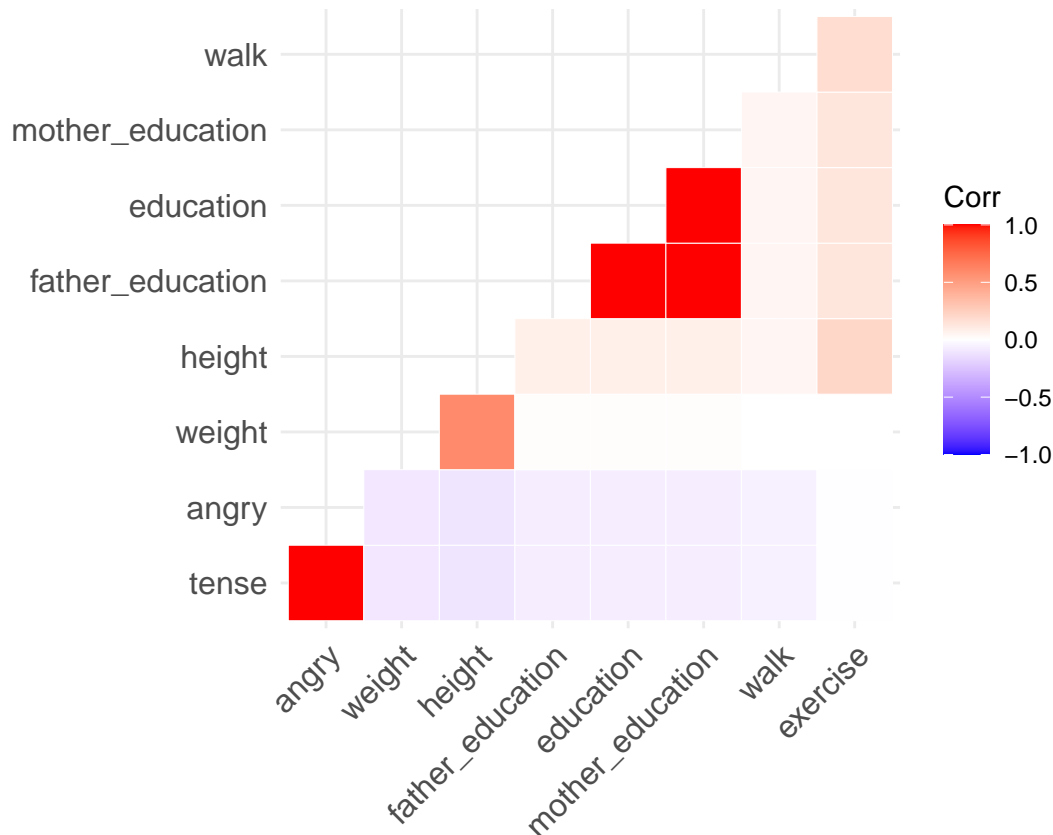
df <- read.csv('earnings.csv')
```

The code below stores a subset of the dataset with the columns indicated above. Since we are calculating the correlation between columns, rows with missing values in any column are removed.

```
subset_cols <- c('education', 'mother_education', 'father_education',
                'walk', 'exercise', 'tense', 'angry', 'weight', 'height')
df_subset <- subset(df, select=subset_cols)
df_subset <- drop_na(df_subset)
```

Created following this guide: [https://rstudio-pubs-static.s3.amazonaws.com/785193\\_5b81a435ca10449ebc16fdf35afd815d.html](https://rstudio-pubs-static.s3.amazonaws.com/785193_5b81a435ca10449ebc16fdf35afd815d.html)

```
corr_matrix <- cor(df_subset)
ggcorrplot(corr_matrix, hc.order=TRUE, type="lower", outline.col = "white")
```



- (10 points) Perform a linear regression model using the variable *earn* as the dependent variable and years of education *education* as the independent variable. What can you say about this covariate? Is it significant? Write down the hypothesis test. Plot the linear regression you have obtained in ggplot by using a subset of the data. This subset is obtained by restricting the variable *earn* to be less than  $2e+05$  (similar to Figure 2 below)
- (20 points) Draw the qqplot by using the library ggplot for the model obtained in point b. Then perform the qqplot (using the library ggplot) for the two different groups of sex (similar to Figure 3 below). Take special care to the legend and the label. What can you say about this plot?
- (20 points) Perform in R the backward and forward procedure to select the covariates, remember to remove the rows with missing values. Did you obtain the same or different results from the two different procedures, please explain. Which procedure would you prefer? Comment what you discovered and the theoretical implications. Just for the backward solution compute the RSS and show the trend of RSS for beta1 in a plot by using ggplot in R (similar to Figure 4). (Hint: For RSS plot, set the range of x-axis to be  $[0,1000]$ ).
- (20 points) Perform a bootstrap of 500 samples for beta 1 (*height*), beta 2 (*male*), and beta 3 (*education*) for the coefficient obtained in the backward procedure in point d. Plot the beta coefficients that you have obtained with histograms with ggplot (similar to Figure 5). Remember to use the data without missing values.
- (20 points) Compute the LOO and K-fold cross validation and write the results. Compute the mean square error for both the LOO and the K-fold cross validation. Then plot the prediction against the true value for LOO, using ggplot. Describe the results. Remember to use the data without missing

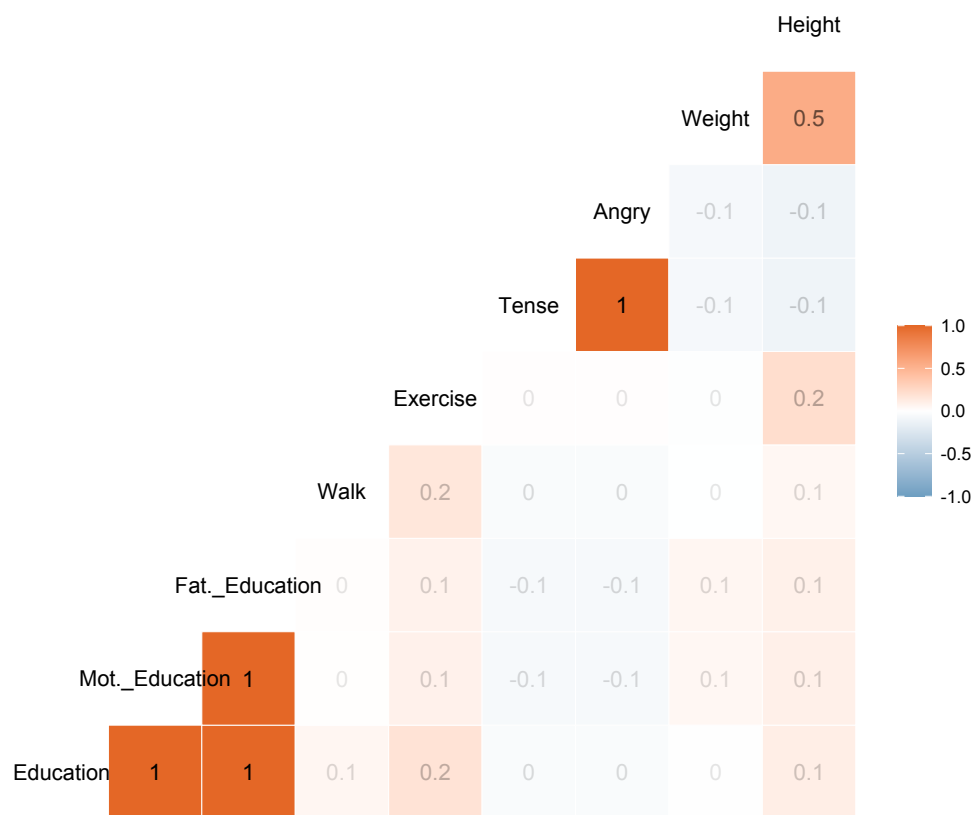


Figure 1: Correlation

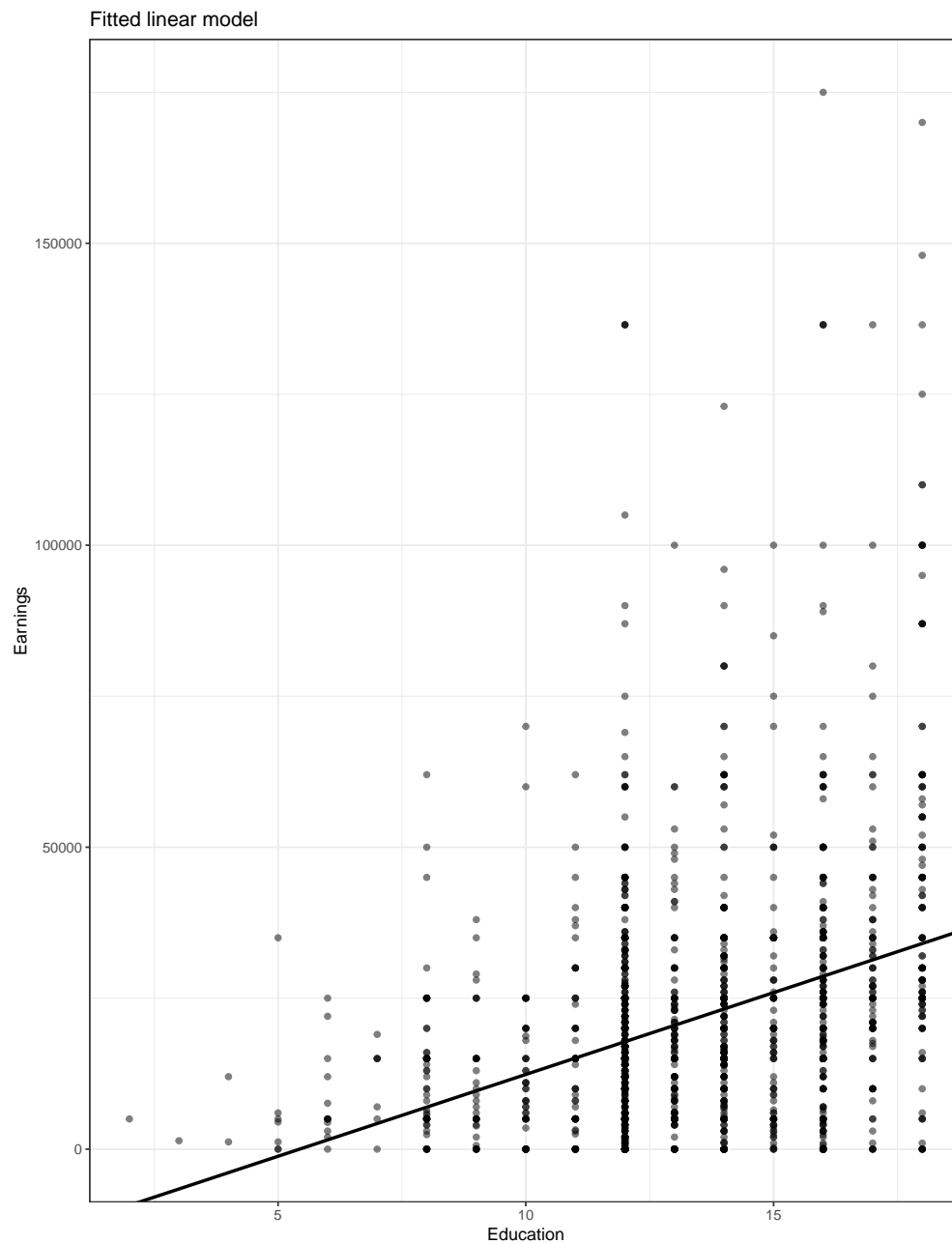


Figure 2: Linear Regression

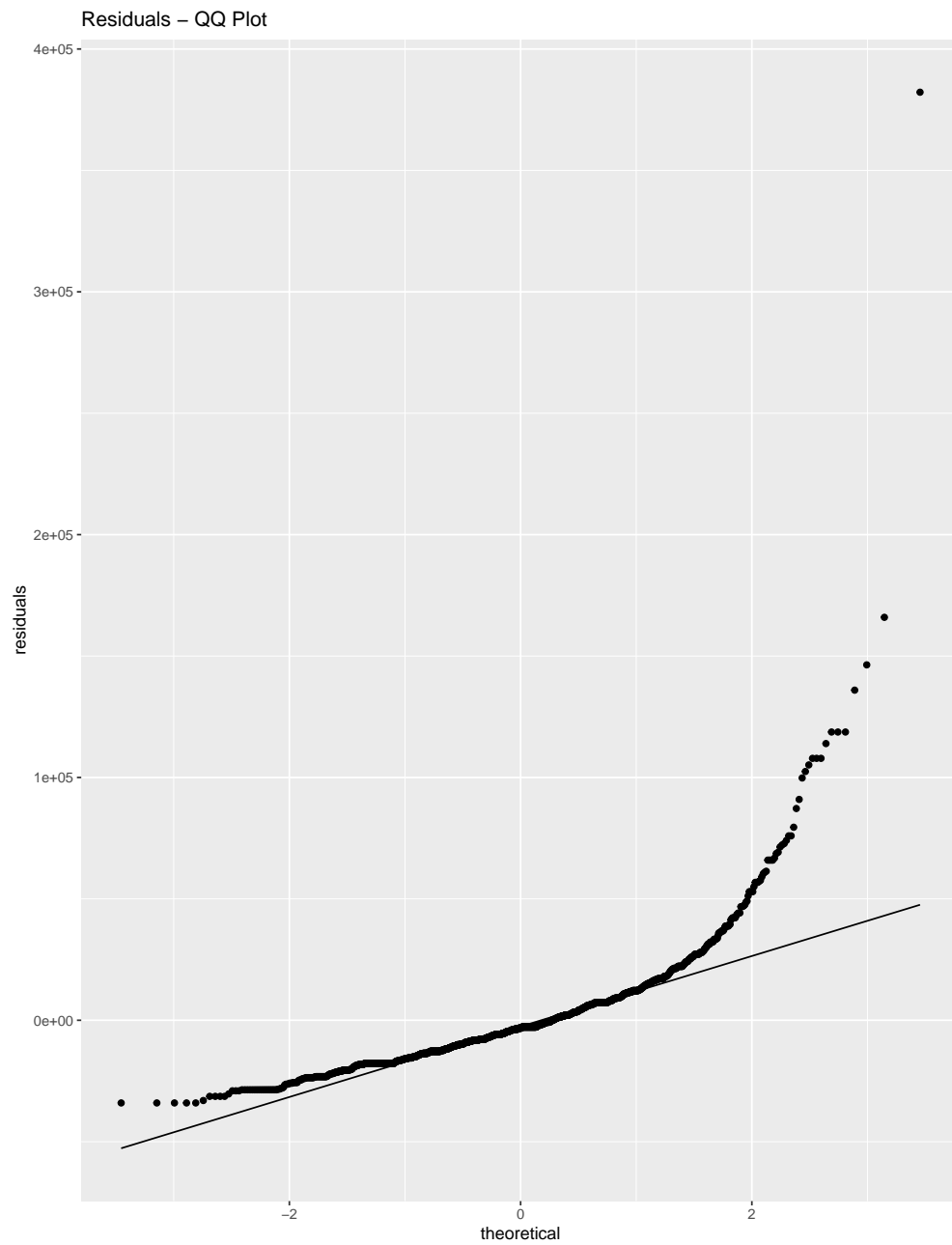


Figure 3: QQplot for different groups

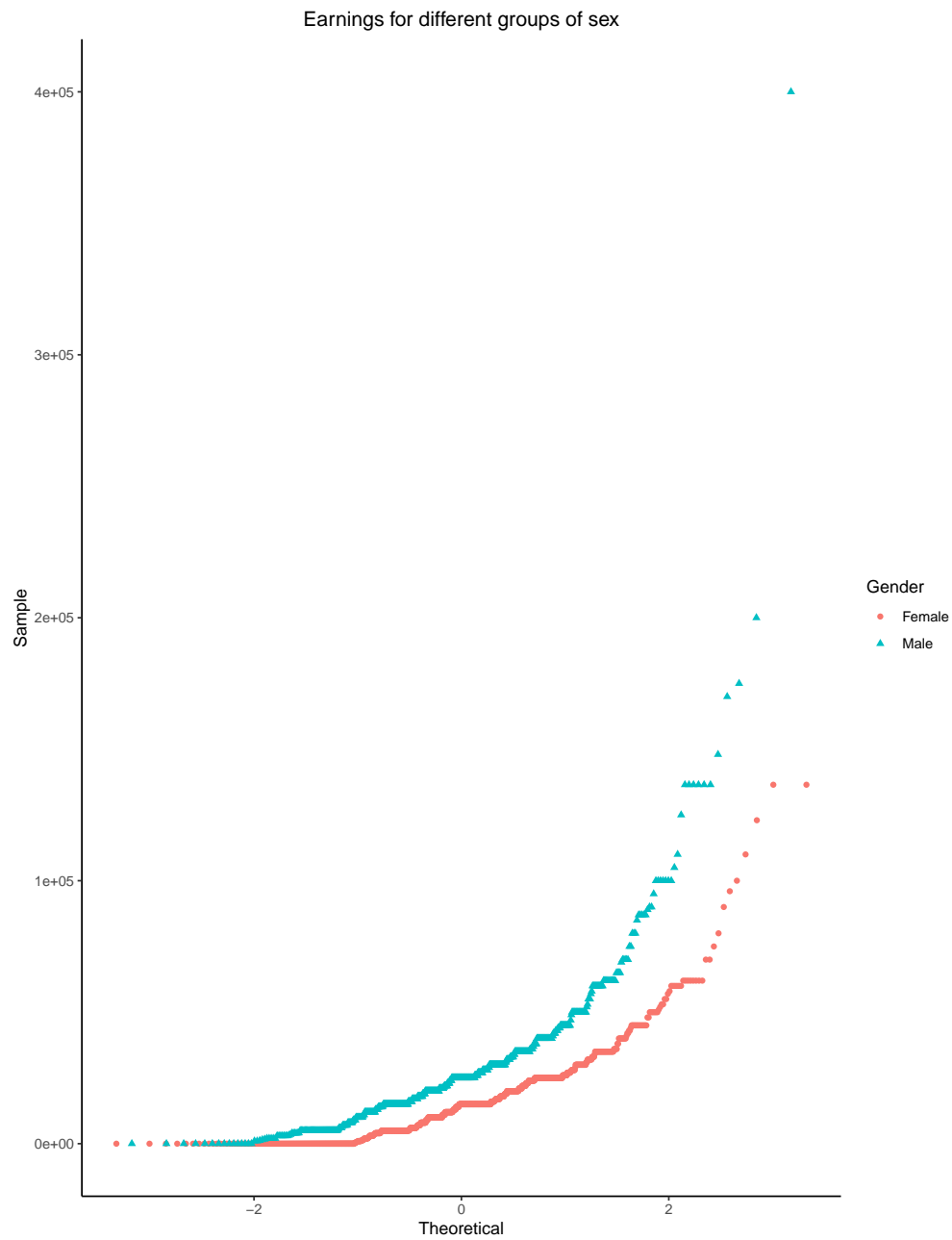


Figure 4: RSS for the backward procedure

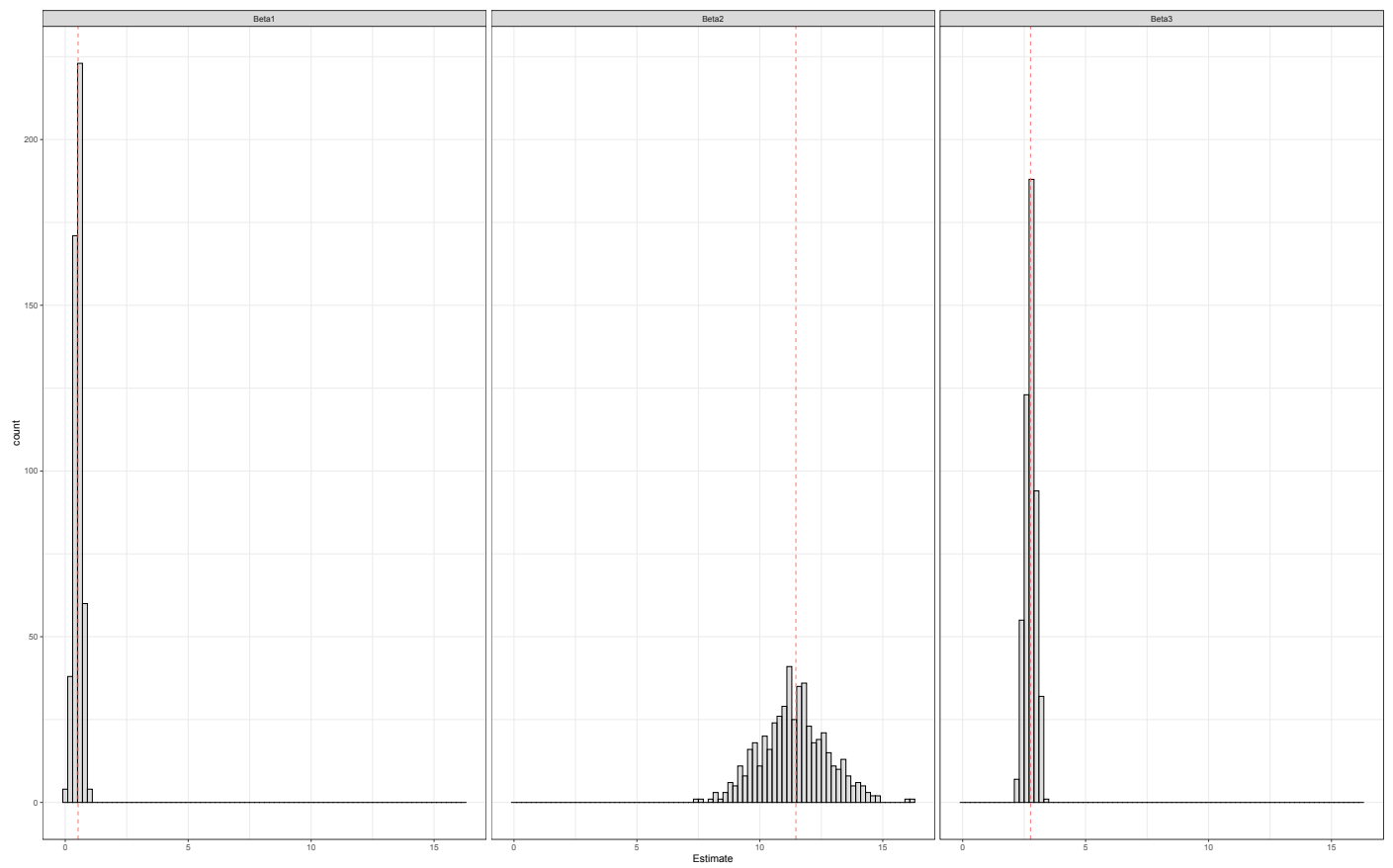


Figure 5: Bootstrap Results

values.