

# Exercise 11 - Linear regression, basic concepts

*Zoltan Kekecs*

*14 november 2019*

## Contents

<b>1</b>	<b>Abstract</b>	<b>2</b>
<b>2</b>	<b>Data management and descriptive statistics</b>	<b>2</b>
2.1	Loading packages . . . . .	2
2.2	Custom functions . . . . .	2
2.3	Get some data . . . . .	3
2.4	Check the dataset for irregularities . . . . .	3
<b>3</b>	<b>Prediction with linear regression</b>	<b>7</b>
3.1	how to set up and interpret simple regression . . . . .	7
3.2	How good is my model . . . . .	10
3.3	Is the model useful? . . . . .	11

# 1 Abstract

This exercise is related to learning the basic logic behind making predictions with linear regression models, and to quantifying the effectiveness of our predictions.

The latest version of this document and the code the document refers to can be found in the GitHub repository of the class at: [https://github.com/kekecsz/PSYP13\\_Data\\_analysis\\_class-2019](https://github.com/kekecsz/PSYP13_Data_analysis_class-2019)

## 2 Data management and descriptive statistics

### 2.1 Loading packages

You will need to load the following packages for this exercise:

```
library(psych)  # for describe
library(gsheet) # to read data from google sheets
library(tidyverse) # for tidy code
```

### 2.2 Custom functions

The following custom functions will come in handy during the visualization of the error of the predictions of the regressions line. This is not essential for the analysis and you don't need to learn how this custom function works if you don't want to.

```
error_plotter <- function(mod, col = "black", x_var = NULL) {
  mod_vars = as.character(mod$call[2])
  data = as.data.frame(eval(parse(text = as.character(mod$call[3]))))
  y = substr(mod_vars, 1, as.numeric(gregexpr(pattern = "~",
    mod_vars)) - 2)
  x = substr(mod_vars, as.numeric(gregexpr(pattern = "~", mod_vars)) +
    2, nchar(mod_vars))

  data$pred = predict(mod)

  if (x == "1" & is.null(x_var)) {
    x = "response_ID"
    data$response_ID = 1:nrow(data)
  } else if (x == "1") {
    x = x_var
  }

  plot(data[, y] ~ data[, x], ylab = y, xlab = x)
  abline(mod)

  for (i in 1:nrow(data)) {
    clip(min(data[, x]), max(data[, x]), min(data[i, c(y,
      "pred")]), max(data[i, c(y, "pred")]))
    abline(v = data[i, x], lty = 2, col = col)
  }
}
```

## 2.3 Get some data

Lets say we are a company that sells shoes, and we would like to be able to tell people's shoe size just by knowing their height. Maybe because our branch is located at a place with lots of tourists, who do not know their european shoe size.

<https://goo.gl/forms/nNXWwCPTbdHxGrBO2>

You can load the data we just collected with the following code

```
mydata = as_tibble(read.csv("https://bit.ly/2Q0qZgP"))
```

## 2.4 Check the dataset for irregularities

You should always check the dataset for coding errors or data that does not make sense.

View data in the data viewer tool

```
View(mydata)
```

and display simple descriptive statistics and plots.

You can investigate the structure of the dataset with the `str()` function, use the `summary()` function to get a basic overview.

```
str(mydata)
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame': 56 obs. of 4 variables:
## $ Timestamp: Factor w/ 56 levels "06/11/2017 16:52:42",...: 1 2 3 4 5 6 7 8 9 10 ...
## $ gender : Factor w/ 5 levels "etc.", "female",...: 4 2 2 2 2 2 2 2 2 ...
## $ height : Factor w/ 32 levels "154", "155", "156",...: 27 9 6 13 19 10 17 2 7 17 ...
## $ shoe_size: int 42 39 36 38 41 38 40 37 39 38 ...
```

```
mydata %>% summary()
```

```
##           Timestamp      gender      height      shoe_size
## 06/11/2017 16:52:42: 1  etc. : 1  167 : 5  Min. :35.00
## 08/11/2017 13:24:57: 1  female:40 172 : 4  1st Qu.:38.00
## 08/11/2017 13:24:59: 1  Female: 2 182 : 4  Median :39.00
## 08/11/2017 13:25:05: 1  male :11 160 : 3  Mean :39.57
## 08/11/2017 13:25:22: 1  Male : 2 165 : 3  3rd Qu.:41.00
## 08/11/2017 13:25:31: 1           155 : 2  Max. :47.00
## (0ther)           :50           (Other):35
```

There are several irregularities in the dataset, such as the height is a factor, and has some weird values, and also the variable 'gender' has inconsistency in coding. We need to handle these before going forward.

You should clean your data before doing the analyses. Make sure to keep a record of the corrections. You should keep the raw datafile when possible, and make changes in the code, so the corrections and omissions are trackable in your code.

```
mydata_cleaned <- mydata %>% mutate(height = as.numeric(as.character(replace(height,
  height == "168cm", 168))), height = as.numeric(as.character(replace(height,
  height == "64", 164))), gender = droplevels(replace(gender,
  gender == "Female", "female")), gender = droplevels(replace(gender,
  gender == "Male", "male")))
```

```
# if the variable gender is identified as a factor, we can
# drop the factor levels that no longer use (such as 'Female'
# and 'Male' instead of 'female' and 'male') using the
# droplevels() function
```

Now we should check whether the data cleaning worked as intended.

You can also use the `describe()` function for more details on the descriptives on each variable (from the `psych` package) which is mostly useful for continuous data, and the `table()` function for categorical data.

```
mydata_cleaned %>% summary()
```

```
##           Timestamp      gender      height      shoe_size
## 06/11/2017 16:52:42: 1   etc.   : 1   Min.   :154.0   Min.   :35.00
## 08/11/2017 13:24:57: 1   female:42   1st Qu.:165.0   1st Qu.:38.00
## 08/11/2017 13:24:59: 1   male  :13   Median :171.0   Median :39.00
## 08/11/2017 13:25:05: 1                      Mean  :171.9   Mean  :39.57
## 08/11/2017 13:25:22: 1                      3rd Qu.:178.5   3rd Qu.:41.00
## 08/11/2017 13:25:31: 1                      Max.   :196.0   Max.   :47.00
## (Other)                :50
```

```
describe(mydata_cleaned)
```

```
##           vars  n   mean    sd median trimmed   mad min max range skew
## Timestamp*    1 56 28.50 16.31   28.5   28.50 20.76    1 56   55 0.00
## gender*        2 56  2.21  0.46    2.0    2.17  0.00    1  3    2 0.77
## height         3 56 171.88 10.24  171.0   171.39  9.64 154 196   42 0.42
## shoe_size      4 56  39.57  2.49   39.0    39.46  2.22  35  47   12 0.62
##           kurtosis   se
## Timestamp*   -1.26 2.18
## gender*       -0.02 0.06
## height        -0.46 1.37
## shoe_size     -0.14 0.33
```

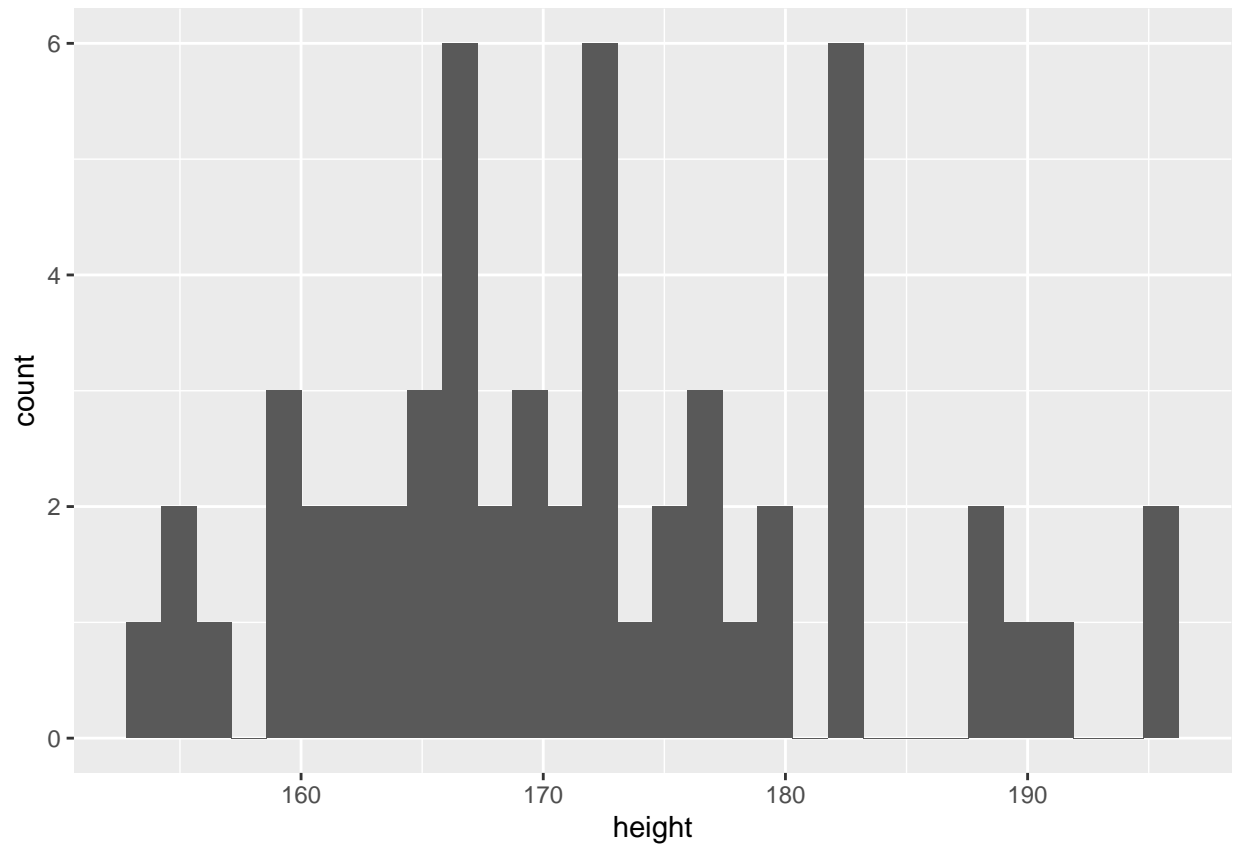
```
table(mydata_cleaned$gender)
```

```
##
##   etc. female   male
##      1     42     13
```

Visualization is also key in checking and exploring data.

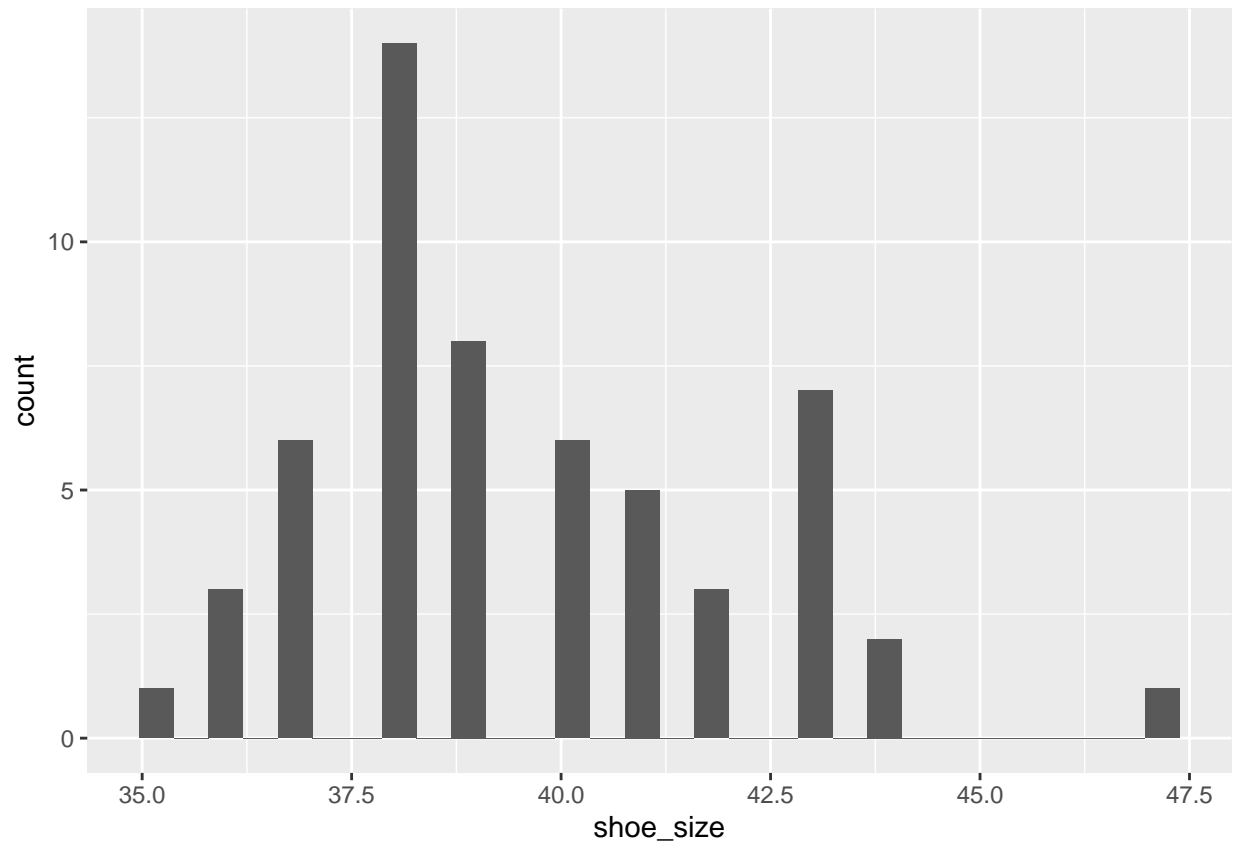
```
mydata_cleaned %>% ggplot() + aes(x = height) + geom_histogram()
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

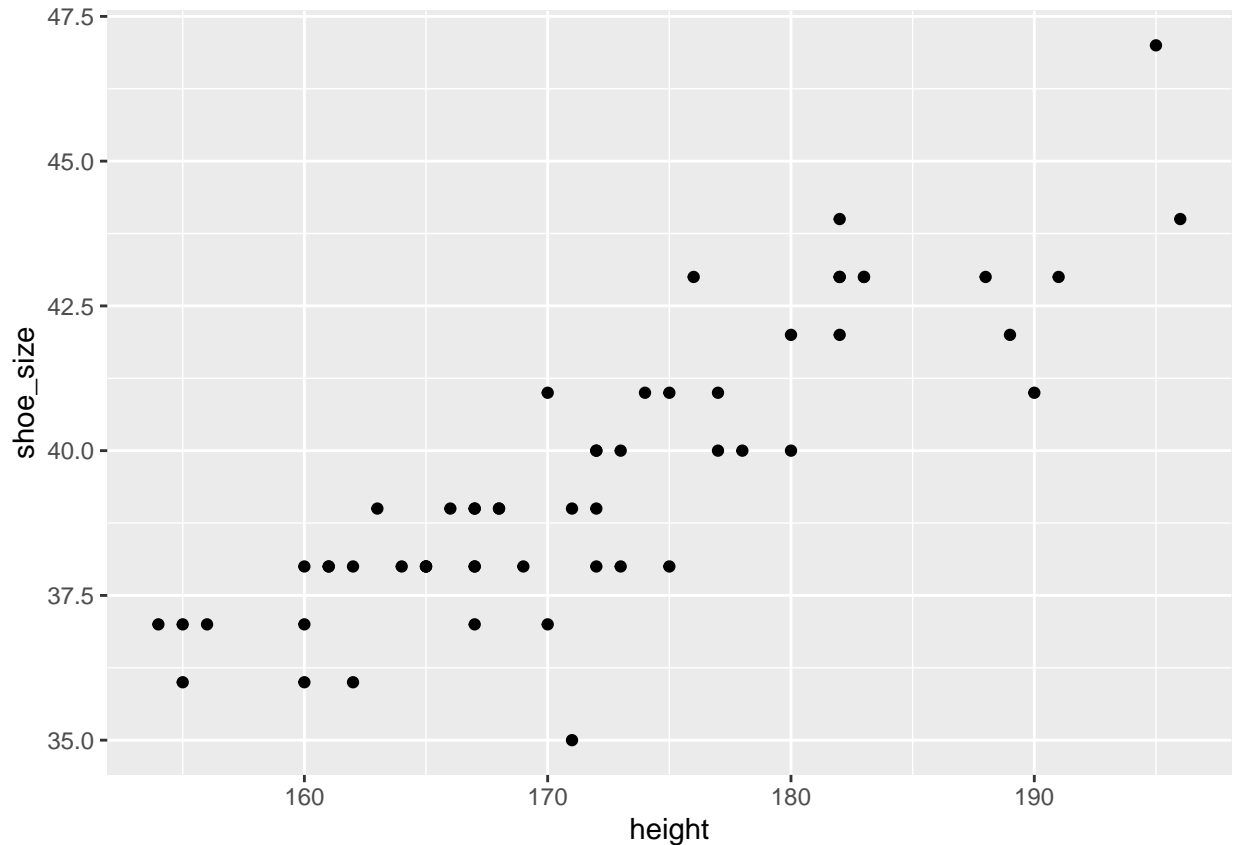


```
mydata_cleaned %>% ggplot() + aes(x = shoe_size) + geom_histogram()

## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



```
mydata_cleaned %>% ggplot() + aes(y = shoe_size, x = height) +  
  geom_point()
```



### 3 Prediction with linear regression

#### 3.1 how to set up and interpret simple regression

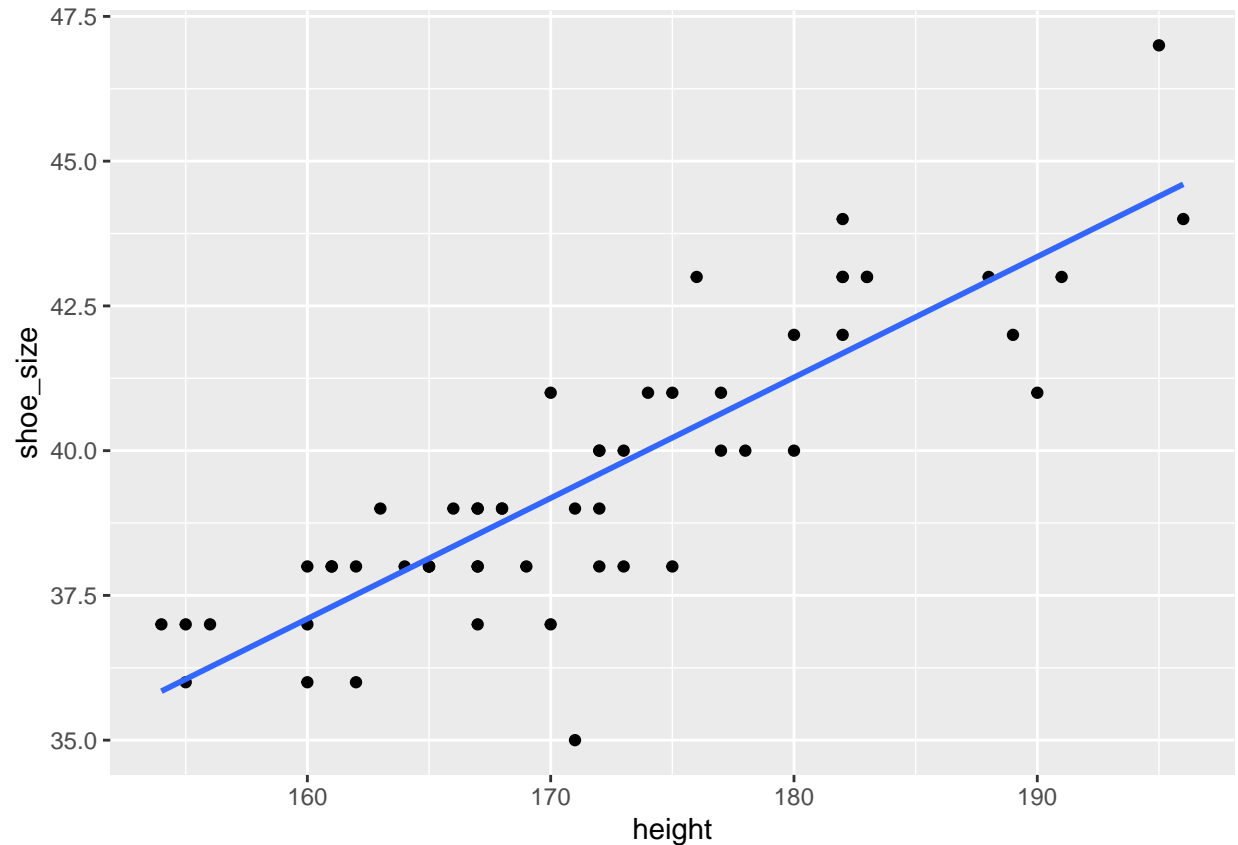
Regression is all about predicting an outcome by knowing the value of predictor variables that are associated with the outcome.

You can set up a simple regression model by using the code below. In the code we first specify an object where the results will be saved (mod1), then we specify the regression formula. In the formula, we start with specifying the outcome variable, what we are interested in predicting (shoe\_size), then, after a ~ we specify the predictors. In simple regression we only have one predictor (here this is height), but we could add more predictor variables, separated with a + sign. In the end of the code we specify the data file where these variables originate from. Remember to use the cleaned dataset (mydata\_cleaned).

```
mod1 <- lm(shoe_size ~ height, data = mydata_cleaned)
```

In simple regression, we identify the underlying pattern of the data, by fitting a single straight line that is closest to all data points.

```
mydata_cleaned %>% ggplot() + aes(x = height, y = shoe_size) +  
  geom_point() + geom_smooth(method = "lm", se = F)
```



Regression provides a mathematical equation (called the regression equation) with which you can predict the outcome by knowing the value of the predictors.

The regression equation is formalized as:  $Y = b_0 + b_1 \cdot X_1$ , where  $Y$  is the predicted value of the outcome,  $b_0$  is the intercept,  $b_1$  is the regression coefficient for predictor 1, and  $X_1$  is the value of predictor 1.

```
mod1

##
## Call:
## lm(formula = shoe_size ~ height, data = mydata_cleaned)
##
## Coefficients:
## (Intercept)      height
##      3.7426      0.2085
```

This means that the regression equation for predicting shoe size is:

shoe size =  $3.74 + 0.21 \cdot \text{height}$

that is, for a person who is 170 cm tall, the predicted shoe size is calculated this way:

$3.74 + 0.21 \cdot 170 = 39.44$

You don't have to do the calculations by hand, you can get the predicted values by using the `predict()` function.

The predictors have to have the same variable name as in the regression formula, and they need to be in a `data.frame` or `tibble` object.



```
height = c(150, 160, 170, 180, 190)
height_df = as_tibble(height)

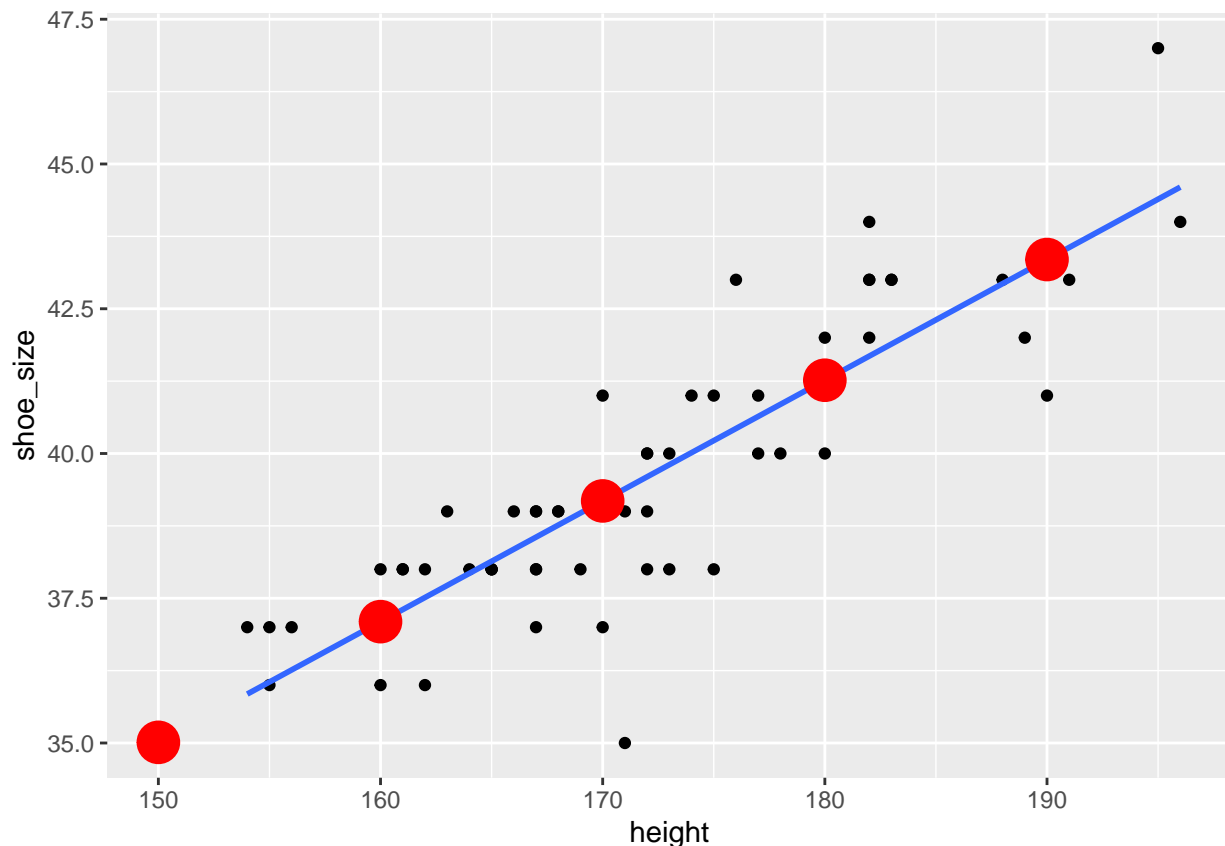
predictions = predict(mod1, newdata = height_df)

height_df_with_predicted = cbind(height_df, predictions)
height_df_with_predicted
```

```
##   value predictions
## 1   150     35.01139
## 2   160     37.09598
## 3   170     39.18057
## 4   180     41.26516
## 5   190     43.34975
```

Predicted values all fall on the regression line

```
mydata_cleaned %>% ggplot() + aes(x = height, y = shoe_size) +
  geom_point() + geom_smooth(method = "lm", se = F) + geom_point(data = height_df_with_predicted,
    aes(x = value, y = predictions), col = "red", size = 7)
```



### Practice

1. Load the mtcars dataset by running `data(mtcars)`. This way you will have a data.frame object called `mtcars` in your workspace. (This data is a built in dataset in R that was extracted from the 1974 Motor Trend US magazine, and comprises fuel consumption and 10 aspects of automobile design and performance for 32 automobiles. You can get more info on it by running `?mtcars`)
2. Build a simple linear regression model using the `lm()` function where **mpg** (miles per gallon) is the

outcome variable and **hp** (raw horse power) is the predictor. Save the output of the model into a new object.

3. Write down the regression equation for estimating mpg
4. Interpret the regression equation. Cars with higher horsepower would have lower or higher miles per gallon? How much higher mpg is predicted for each step up on hp?
5. Create a scatterplot of the relationship of mpg and hp including the regression line and use the plot to verify your interpretation from before.
6. Using this model, what is the predicted mpg for cars with horsepower of 100, 200, and 300? (you can use the `predict()` function to answer this question)

---

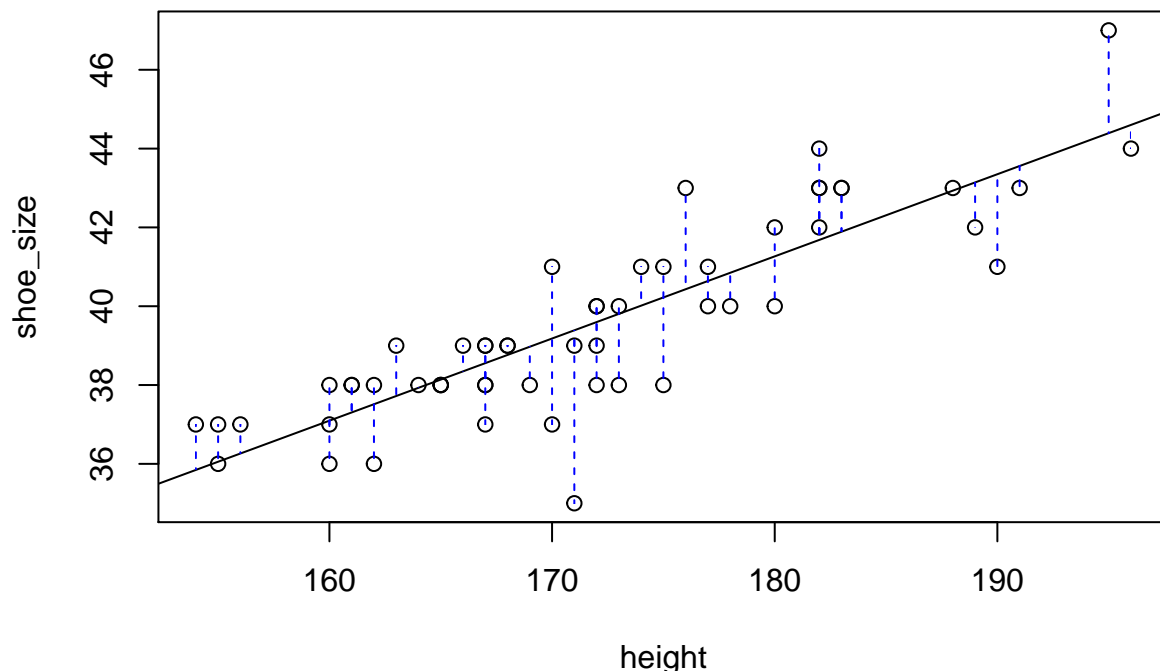
## 3.2 How good is my model

### 3.2.1 How to measure prediction efficiency?

You can measure how effective your model is by measuring the difference between the actual outcome values and the predicted values. We call this residual error in regression.

The residual error for each observed shoe size can be seen on this plot (this will only work if you ran the code in the top of the script containing the `error_plotter()` custom function). This code is only for visualization purposes.

```
error_plotter(mod1, col = "blue")
```



You can simply add up all the residual error (the length of these lines), and get a good measure of the overall efficiency of your model. This is called the residual absolute difference (RAD). However, this value is rarely used. More common is the residual sum of squared differences (RSS). The interpretation is practically the same, gives an indication of the total amount of error when using the model.

```
RAD = sum(abs(mydata_cleaned$shoe_size - predict(mod1)))
RAD
```

```
## [1] 54.79272
```

```
RSS = sum((mydata_cleaned$shoe_size - predict(mod1))^2)
RSS
```

```
## [1] 91.1467
```

### 3.3 Is the model useful?

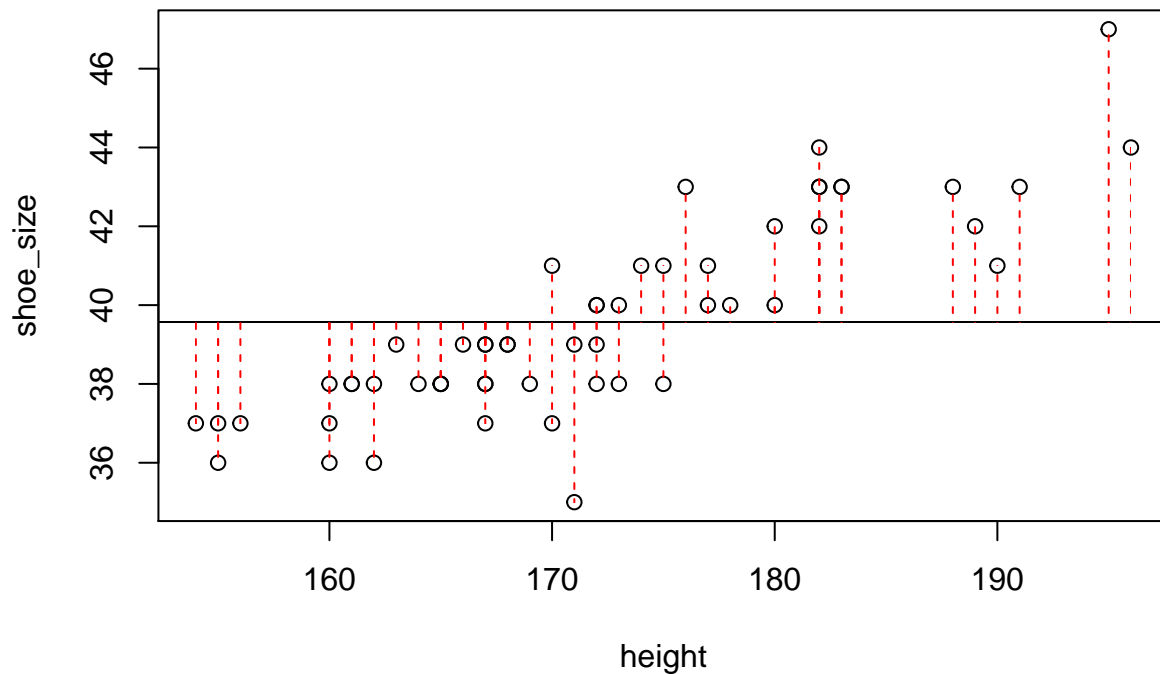
To establish how much benefit did we get by taking into account the predictor, we can compare the residual error when using our best guess (mean) without taking into account the predictor, with the residual error when the predictor is taken into account.

Below you can find regression model where we only use the mean of the outcome variable to predict the outcome value.

We can calculate the sum of squared differences the same way as before, but for the model without any predictors, we call this the total sum of squared differences (TSS).

```
mod_mean <- lm(shoe_size ~ 1, data = mydata_cleaned)
```

```
error_plotter(mod_mean, col = "red", x_var = "height") # visualize error
```



```
TSS = sum((mydata_cleaned$shoe_size - predict(mod_mean))^2)
TSS
```

```
## [1] 341.7143
```

The total amount of information gained about the variability of the outcome is shown by the R squared statistic ( $R^2$ ).

```
R2 = 1 - (RSS/TSS)
R2
```

```
## [1] 0.7332663
```

This means that by using the regression model, we are able to explain 73.33% of the variability in the outcome.

$R^2 = 1$  means all variability of the outcome is perfectly predicted by the predictor(s)

$R^2 = 0$  means no variability of the outcome is predicted by the predictor(s)

### 3.3.1 Is the model with the predictor significantly better than a model without the predictor?

You can do an anova to find out, comparing the amount of variance explained by the two models.

```
anova(mod_mean, mod1)
```

```
## Analysis of Variance Table
##
## Model 1: shoe_size ~ 1
## Model 2: shoe_size ~ height
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      55 341.71
## 2      54  91.15  1    250.57 148.45 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Or, you can get all this information and more from the model summary

```
summary(mod1)
```

```
##
## Call:
## lm(formula = shoe_size ~ height, data = mydata_cleaned)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.3890 -0.6103  0.2152  0.7477  2.6080
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3.74256    2.94578   1.27    0.209
## height        0.20846    0.01711  12.18 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.299 on 54 degrees of freedom
## Multiple R-squared:  0.7333, Adjusted R-squared:  0.7283
## F-statistic: 148.4 on 1 and 54 DF,  p-value: < 2.2e-16
```

The confidence interval of the regression coefficient is given by the confint() function

```
confint(mod1)
```

```
##              2.5 %      97.5 %
```

```
## (Intercept) -2.1633697 9.6484844  
## height      0.1741569 0.2427609
```

You can also plot the confidence interval of the predictions using `geom_smooth()` in `ggplot()`

```
mydata_cleaned %>% ggplot() + aes(x = height, y = shoe_size) +  
  geom_point() + geom_smooth(method = "lm")
```

