

# Simple Linear Regression

Date: 2017-11-28

## What it is

In its simplest form, simple linear regression is used to predict the value of a numerical **outcome**<sup>1</sup> ( $y$ ) from a numerical **predictor**<sup>2</sup> ( $x$ ). The predicted value  $\hat{y}$  will have the smallest squared error  $(y_i - \hat{y})^2$  of any linear predictor based on only  $x$ .

<sup>1</sup> Also called the “response” or “dependent variable”

<sup>2</sup> Also called the “explanatory variable” or “independent variable”

## When to use it

Linear regression is used to estimate the functional form of a linear relationship between two variables<sup>3</sup>. In other words, it estimates the slope and intercept of the line that best represents the trend between the variables. As with most statistical methods, we make some simplifying assumptions when using simple linear regression which should be checked before the regression model is put forth as a potential solution<sup>4</sup>. These assumptions are:

<sup>3</sup> As opposed to correlation, which estimates the *strength* of the relationship

<sup>4</sup> See *Regression Diagnostics*

1. There is a linear trend in the data
2. Errors/residuals are normally distributed
3. Errors/residuals have a constant variance across the range of  $x$  and  $y$ .

## How to use it

To run a linear regression in R, use the `lm()` function<sup>5</sup>. The regression variables should be included in a formula of the form  $y \sim x$ , along with the data set. The regression model is then saved for future use. For example, if we want to use height to predict weight, we would use

<sup>5</sup> `lm` stands for “linear model.”

```
myreg = lm(weight ~ height, data = d)
```

Most information you would want from the regression model can be obtained through the `summary()` function. This output will include intercept and slope calculations with hypothesis tests for each, the residual standard deviation,  $R^2$ , and an  $F$  test for model fit:

```
summary(myreg)
```

```
##
```

```
## Call:
```

```
## lm(formula = weight ~ height, data = d)
```

```
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -107.27  -28.46   -6.90   18.31   324.95
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -180.845     31.426  -5.755 1.52e-08 ***
## height       5.351       0.468  11.435 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 42.51 on 498 degrees of freedom
## Multiple R-squared:  0.208, Adjusted R-squared:  0.2064
## F-statistic: 130.7 on 1 and 498 DF, p-value: < 2.2e-16
```

The output above tells us the estimated regression model is

$$weight = -180.845 + 5.351(height)$$

and the residuals have a standard deviation of 42.51.

*Why it works*