

A Perspective on Multiple Regression

Graduate Statistics

February 19, 2016

Consider a typical multiple regression equation:

$$Y_i = \beta_0 + \beta_x X_i + \beta_z Z_i + \varepsilon_i \quad (1)$$

We interpret the coefficient β_1 as representing the *independent contributions of X to Y* or as *what X uniquely explains in Y independent of Z*. Here we attempt this interpretation precise.

Multiple regression in terms of residuals.

A slope in multiple regression represents the relationship between *a*) the portion of a predictor not explained by the other predictors and *b*) the portion of the outcome variable not explained by those same other predictors. The coefficient β_z in equation (1) can thus reproduced as follows:

First, we regress Z on X :

$$Z_i = \alpha_0 + \alpha_1 X_i + \varepsilon_{zx_i} \quad (2)$$

We also regress Y on X :

$$Y_i = \alpha_0 + \alpha_1 X_i + \varepsilon_{yx_i} \quad (3)$$

Now we can examine our residual error terms. ε_{zx} is the portion of Z that wasn't explained by X . That is, we estimated model (2) that was able to explain some of the variability in Z , but not all of it; ε_{zx} is what's left over. Analogously, ε_{yx} is the portion of Y that wasn't explained by X . Now let's model the relationship between these two sets of residuals¹:

$$\varepsilon_{yx} = \alpha_0 + \alpha_z \varepsilon_{zx} + \gamma_i \quad (4)$$

Here, In the context of our full model (1), we would say that β_z represents the relationship between the part of Z not explained by X and the part of Y not explained by X . In other words, β_z in model (1) is the same α_z in model (4). Let's look at a concrete example of this.

Horsepower, engine size, and fuel efficiency.

For the following example, we will use the `mtcars` dataset that comes with R. Here is the description from the help page:

The data was extracted from the 1974 Motor Trend US magazine, and comprises fuel consumption and 10 aspects of automobile design and performance for 32 automobiles (1973-74 models).

Here, we will be regressing miles per gallon (`mpg`) on horsepower (`hp`) and engine size (`disp`):

$$\text{mpg}_i = \beta_0 + \beta_d \text{disp}_i + \beta_h \text{hp}_i + \gamma_i \quad (5)$$

¹ γ (pronounced "gamma") here is just another error term. A different letter is used to differentiate the residuals from the predictor.

```
fullModel <- lm(mpg ~ disp + hp, data = mtcars)
summary(fullModel)$coefficients

##           Estimate Std. Error  t value    Pr(>|t|)
## (Intercept) 30.73590425 1.331566129 23.082522 3.262507e-20
## disp      -0.03034628 0.007404856 -4.098159 3.062678e-04
## hp        -0.02484008 0.013385499 -1.855746 7.367905e-02
```

Now we will verify that equation (4) holds. We'll extract the residuals from each of the intermediate models:

$$\text{mpg}_i = \alpha_0 + \alpha_1 \text{hp}_i + \varepsilon_{mh_i} \quad (6)$$

$$\text{disp}_i = \alpha_0 + \alpha_1 \text{hp}_i + \varepsilon_{dh_i} \quad (7)$$

```
mpgONhp <- lm(mpg ~ hp, data = mtcars)$residuals
dispONhp <- lm(disp ~ hp, data = mtcars)$residuals
```

We now confirm that β_d in equation (5) represents the relationship between the error terms in equations (6) and (7)²:

$$\varepsilon_{mh_i} = \beta_0 + \beta_d \varepsilon_{dh_i} + \gamma_i$$

```
equivalent <- lm(mpgONhp ~ dispONhp)
summary(equivalent)$coefficients

##           Estimate Std. Error  t value    Pr(>|t|)
## (Intercept) 3.879825e-16 0.543420336 7.139640e-16 1.0000000000
## dispONhp   -3.034628e-02 0.007280396 -4.168218e+00 0.0002400329

summary(fullModel)$coefficients

##           Estimate Std. Error  t value    Pr(>|t|)
## (Intercept) 30.73590425 1.331566129 23.082522 3.262507e-20
## disp      -0.03034628 0.007404856 -4.098159 3.062678e-04
## hp        -0.02484008 0.013385499 -1.855746 7.367905e-02
```

As expected, the coefficients are the same. The test statistic values only differ because the standard errors of the estimates depend on the remaining degrees of freedom for the full model (the t -value for our equivalent model is a little higher than it should be because we're ignoring the fact that we used up an extra degree of freedom when we ran models (6) and (7).

Note: this document was created with knitr and LyX.

²Note that the intercept here will always be zero as residuals always have a mean of zero.