

# STA221

Neil Montgomery

Last edited: 2017-03-22 15:05

multiple regression

## regression with more than one input variable

The Universal Statistical Model:

$$\text{Output} = \text{Input} + \text{Noise}$$

## regression with more than one input variable

The Universal Statistical Model:

$$\text{Output} = \text{Input} + \text{Noise}$$

Most datasets have more than one or two columns.

## regression with more than one input variable

The Universal Statistical Model:

$$\text{Output} = \text{Input} + \text{Noise}$$

Most datasets have more than one or two columns.

The most important statistical model (in my opinion) is the linear regression model with more than one “x” variable. For example, with 3 input variables:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \varepsilon$$

## interpretation of the variables

We treat  $y$  as random. The inputs are not random. They can be whatever you like, even functions of one another, with one technical limitation\*.

## interpretation of the variables

We treat  $y$  as random. The inputs are not random. They can be whatever you like, even functions of one another, with one technical limitation\*.

So, for example, the following is a valid multiple regression model:

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \varepsilon$$

This kind of “polynomial” model is good for fitting some types of non-linear relationships between  $y$  and a single  $x$ .

## interpretation of the variables

We treat  $y$  as random. The inputs are not random. They can be whatever you like, even functions of one another, with one technical limitation\*.

So, for example, the following is a valid multiple regression model:

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \varepsilon$$

This kind of “polynomial” model is good for fitting some types of non-linear relationships between  $y$  and a single  $x$ .

\*A variable cannot be a linear function of other variables in the model.



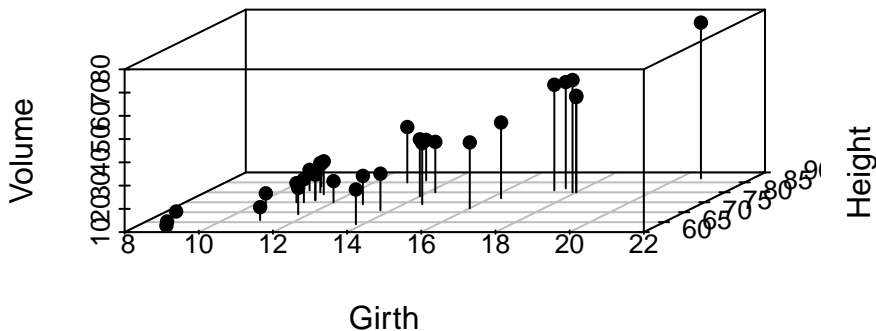
## what is being accomplished in multiple regression?

R comes with some sample datasets. One is called `trees` and has variables `Girth`, `Height`, and `Volume`. Here's a peek at the data:

```
## # A tibble: 31 × 3
##   Girth Height Volume
##   <dbl>  <dbl>  <dbl>
## 1    8.3     70   10.3
## 2    8.6     65   10.3
## 3    8.8     63   10.2
## 4   10.5     72   16.4
## 5   10.7     81   18.8
## # ... with 26 more rows
```

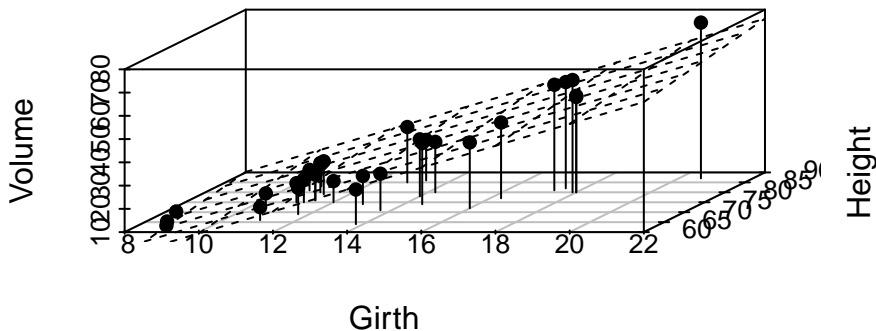
what is being accomplished in multiple regression?

## Volume versus height and girth



multiple regression fits a surface to the points

## Volume versus height and girth



## the fundamental issues

- ▶ Familiar issues with similar answers

## the fundamental issues

- ▶ Familiar issues with similar answers
  - ▶ Parameter testing and estimation

## the fundamental issues

- ▶ Familiar issues with similar answers
  - ▶ Parameter testing and estimation
  - ▶ Mean response and prediction

## the fundamental issues

- ▶ Familiar issues with similar answers
  - ▶ Parameter testing and estimation
  - ▶ Mean response and prediction
  - ▶ Model assumptions

## the fundamental issues

- ▶ Familiar issues with similar answers
  - ▶ Parameter testing and estimation
  - ▶ Mean response and prediction
  - ▶ Model assumptions
- ▶ New issues:



## the fundamental issues

- ▶ Familiar issues with similar answers
  - ▶ Parameter testing and estimation
  - ▶ Mean response and prediction
  - ▶ Model assumptions
- ▶ New issues:
  - ▶ Parameter interpretation

## the fundamental issues

- ▶ Familiar issues with similar answers
  - ▶ Parameter testing and estimation
  - ▶ Mean response and prediction
  - ▶ Model assumptions
- ▶ New issues:
  - ▶ Parameter interpretation
  - ▶ Hard to visualize what is really happening

## the fundamental issues

- ▶ Familiar issues with similar answers
  - ▶ Parameter testing and estimation
  - ▶ Mean response and prediction
  - ▶ Model assumptions
- ▶ New issues:
  - ▶ Parameter interpretation
  - ▶ Hard to visualize what is really happening
  - ▶ Actual formulae too unwieldy to even present

## the fundamental issues

- ▶ Familiar issues with similar answers
  - ▶ Parameter testing and estimation
  - ▶ Mean response and prediction
  - ▶ Model assumptions
- ▶ New issues:
  - ▶ Parameter interpretation
  - ▶ Hard to visualize what is really happening
  - ▶ Actual formulae too unwieldy to even present
  - ▶ Model selection: which variables?

## the fundamental issues

- ▶ Familiar issues with similar answers
  - ▶ Parameter testing and estimation
  - ▶ Mean response and prediction
  - ▶ Model assumptions
- ▶ New issues:
  - ▶ Parameter interpretation
  - ▶ Hard to visualize what is really happening
  - ▶ Actual formulae too unwieldy to even present
  - ▶ Model selection: which variables?
  - ▶ “Multicollinearity” (highly correlated inputs)

## parameter interpretation

The multiple regression model:

$$y = \beta_0 + \beta_1 x_1 + \dots \beta_k x_k + \varepsilon, \quad \varepsilon \sim N(0, \sigma)$$

has many parameters.

## parameter interpretation

The multiple regression model:

$$y = \beta_0 + \beta_1 x_1 + \dots \beta_k x_k + \varepsilon, \quad \varepsilon \sim N(0, \sigma)$$

has many parameters.

$\sigma$  is the variation in the distribution of the noise. It is not a function of any of the  $x$  - just like before it is a constant.

## parameter interpretation

The multiple regression model:

$$y = \beta_0 + \beta_1 x_1 + \dots \beta_k x_k + \varepsilon, \quad \varepsilon \sim N(0, \sigma)$$

has many parameters.

$\sigma$  is the variation in the distribution of the noise. It is not a function of any of the  $x$  - just like before it is a constant.

$\beta_0$  is the “intercept”—mainly important to make sure the fitted surface actually goes through the points.



## parameter interpretation

The multiple regression model:

$$y = \beta_0 + \beta_1 x_1 + \dots \beta_k x_k + \varepsilon, \quad \varepsilon \sim N(0, \sigma)$$

has many parameters.

$\sigma$  is the variation in the distribution of the noise. It is not a function of any of the  $x$  - just like before it is a constant.

$\beta_0$  is the “intercept”—mainly important to make sure the fitted surface actually goes through the points.

The  $\beta_i$  from  $i \in \{1, \dots, k\}$  are the slope parameters, and have a different interpretation than before.

## slope parameter interpretation

$\beta_i$  is:

- ▶ the change in  $y$

## slope parameter interpretation

$\beta_i$  is:

- ▶ the change in  $y$
- ▶ when  $x_i$  increases by 1 unit

## slope parameter interpretation

$\beta_i$  is:

- ▶ the change in  $y$
- ▶ when  $x_i$  increases by 1 unit
- ▶ ***given [values of] all the other input variables in the model.***

## slope parameter interpretation

$\beta_i$  is:

- ▶ the change in  $y$
- ▶ when  $x_i$  increases by 1 unit
- ▶ ***given [values of] all the other input variables in the model.***

## slope parameter interpretation

$\beta_i$  is:

- ▶ the change in  $y$
- ▶ when  $x_i$  increases by 1 unit
- ▶ ***given [values of] all the other input variables in the model.***

That bold, italic statement should echo in your mind any time you think of anything to do with  $\beta_i$ .

## trees example

We might want to model  $y = \text{Volume}$  (the amount of wood) as a linear model of the input variables  $x_1 = \text{Girth}$  and  $x_2 = \text{Height}$ , as follows:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$$

## trees example

We might want to model  $y = \text{Volume}$  (the amount of wood) as a linear model of the input variables  $x_1 = \text{Girth}$  and  $x_2 = \text{Height}$ , as follows:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$$

The computer does all the estimation of the parameters.



## trees example

We might want to model  $y = \text{Volume}$  (the amount of wood) as a linear model of the input variables  $x_1 = \text{Girth}$  and  $x_2 = \text{Height}$ , as follows:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$$

The computer does all the estimation of the parameters.

We'll call the fitted model:

$$y = b_0 + b_1 x_1 + b_2 x_2$$

## trees example

We might want to model  $y = \text{Volume}$  (the amount of wood) as a linear model of the input variables  $x_1 = \text{Girth}$  and  $x_2 = \text{Height}$ , as follows:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$$

The computer does all the estimation of the parameters.

We'll call the fitted model:

$$y = b_0 + b_1 x_1 + b_2 x_2$$

The computer uses the method of “least squares”, like before. A full treatment of the analysis requires matrix algebra.

## fitted values | residuals

Here's the first row of the trees data:

Girth	Height	Volume
8.3	70	10.3

We could call these values  $y_1$ ,  $x_{11}$ , and  $x_{21}$

## fitted values | residuals

Here's the first row of the trees data:

Girth	Height	Volume
8.3	70	10.3

We could call these values  $y_1$ ,  $x_{11}$ , and  $x_{21}$

The fitted value for  $y_1$  is just:

$$\hat{y}_1 = b_0 + b_1 x_{11} + b_2 x_{21}$$

## fitted values | residuals

Here's the first row of the trees data:

Girth	Height	Volume
8.3	70	10.3

We could call these values  $y_1$ ,  $x_{11}$ , and  $x_{21}$

The fitted value for  $y_1$  is just:

$$\hat{y}_1 = b_0 + b_1 x_{11} + b_2 x_{21}$$

The residual corresponding to this fitted value is just:

$$y_1 - \hat{y}_1$$

## fitted values | residuals

Here's the first row of the trees data:

Girth	Height	Volume
8.3	70	10.3

We could call these values  $y_1$ ,  $x_{11}$ , and  $x_{21}$

The fitted value for  $y_1$  is just:

$$\hat{y}_1 = b_0 + b_1 x_{11} + b_2 x_{21}$$

The residual corresponding to this fitted value is just:

$$y_1 - \hat{y}_1$$

For a dataset with  $n$  rows (the sample size), there is a fitted value and residual for each row.

## trees data fitted model

Here's what R produces:

```
##  
## Coefficients:  
##           Estimate Std. Error t value Pr(>|t|)  
## (Intercept) -57.9877      8.6382  -6.713 2.75e-07  
## Girth        4.7082       0.2643  17.816 < 2e-16  
## Height       0.3393       0.1302   2.607  0.0145  
##  
## Residual standard error: 3.882 on 28 degrees of freedom  
## Multiple R-squared:  0.948, Adjusted R-squared:  0.9442  
## F-statistic: 255 on 2 and 28 DF, p-value: < 2.2e-16
```

## individual slope parameter hypothesis testing

The usual hypothesis test for a single parameter:

$$H_0 : \beta_i = 0$$

$$H_a : \beta_i \neq 0$$



## individual slope parameter hypothesis testing

The usual hypothesis test for a single parameter:

$$H_0 : \beta_i = 0$$

$$H_a : \beta_i \neq 0$$

If  $H_0$  is true, it means the  $i$ th variable ( $x_i$ ) is not significantly related to  $y$

## individual slope parameter hypothesis testing

The usual hypothesis test for a single parameter:

$$H_0 : \beta_i = 0$$

$$H_a : \beta_i \neq 0$$

If  $H_0$  is true, it means the  $i$ th variable ( $x_i$ ) is not significantly related to  $y$   
***given all the other  $x$ 's in the model***

## the overall hypothesis test

“Is there any linear relationship between  $y$  and the input variables?”

## the overall hypothesis test

“Is there any linear relationship between  $y$  and the input variables?”

Null hypothesis can be expressed as:

$$H_0 : \beta_1 = \beta_2 = \cdots = \beta_k = 0$$

## the overall hypothesis test

“Is there any linear relationship between  $y$  and the input variables?”

Null hypothesis can be expressed as:

$$H_0 : \beta_1 = \beta_2 = \cdots = \beta_k = 0$$

It is also possible to test any subset of these parameters, such as:

$$H_0 : \beta_1 = \beta_2 = 0$$

although at the moment it's not clear why this might be a good idea.

## estimating $\sigma$

This works the same as with simple regression, in which we used  $\sqrt{MSE}$  where:

$$MSE = \frac{\sum_{j=1}^n (y_j - \hat{y}_j)^2}{n - 2}$$

## estimating $\sigma$

This works the same as with simple regression, in which we used  $\sqrt{MSE}$  where:

$$MSE = \frac{\sum_{j=1}^n (y_j - \hat{y}_j)^2}{n - 2}$$

$n - 2$  was the sample size minus the number of parameters (two:  $\beta_0$  and  $\beta_1$ ) being estimated.

## estimating $\sigma$

This works the same as with simple regression, in which we used  $\sqrt{MSE}$  where:

$$MSE = \frac{\sum_{j=1}^n (y_j - \hat{y}_j)^2}{n - 2}$$

$n - 2$  was the sample size minus the number of parameters (two:  $\beta_0$  and  $\beta_1$ ) being estimated.

There was only one input variable, so another way to think of this was “sample size minus the number of input variables, then minus 1.”



## estimating $\sigma$

In multiple regression, nothing changes. Use  $\sqrt{MSE}$ , where:

$$MSE = \frac{\sum_{j=1}^n (y_j - \hat{y}_j)^2}{n - (k + 1)}$$

## hypothesis testing for $\beta_i$

The computer produces the estimate  $b_i$ , which has these properties:

$$E(b_i) = \beta_i$$

$$\text{Var}(b_i) = \sigma \cdot c_i$$

## hypothesis testing for $\beta_i$

The computer produces the estimate  $b_i$ , which has these properties:

$$E(b_i) = \beta_i$$

$$\text{Var}(b_i) = \sigma \cdot c_i$$

$c_i$  is a number that reflects the relationships between  $x_i$  and the other inputs (to be revisited).

## hypothesis testing for $\beta_i$

The computer produces the estimate  $b_i$ , which has these properties:

$$E(b_i) = \beta_i$$

$$\text{Var}(b_i) = \sigma \cdot c_i$$

$c_i$  is a number that reflects the relationships between  $x_i$  and the other inputs (to be revisited).

Just like before, we get:

$$\frac{b_i - \beta_i}{\sqrt{MSE} \sqrt{c_i}} \sim t_{n-k+1}$$

## hypothesis testing for $\beta_i$ in the trees example

```
##  
## Coefficients:  
##           Estimate Std. Error t value Pr(>|t|)  
## (Intercept) -57.9877      8.6382  -6.713 2.75e-07  
## Girth        4.7082       0.2643  17.816 < 2e-16  
## Height       0.3393       0.1302   2.607  0.0145  
##  
## Residual standard error: 3.882 on 28 degrees of freedom  
## Multiple R-squared:  0.948, Adjusted R-squared:  0.9442  
## F-statistic: 255 on 2 and 28 DF, p-value: < 2.2e-16
```