

# Multiple Linear Regression

When we want to understand the relationship between a single predictor variable and a response variable, we often use simple linear regression.

However, if we'd like to understand the relationship between *multiple* predictor variables and a response variable then we can instead use **multiple linear regression**.

If we have  $p$  predictor variables, then a multiple linear regression model takes the form:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \epsilon$$

where:

- $Y$ : The response variable
- $X_j$ : The  $j$ th predictor variable
- $\beta_j$ : The average effect on  $Y$  of a one unit increase in  $X_j$ , holding all other predictors fixed
- $\epsilon$ : The error term

The values for  $\beta_0, \beta_1, \beta_2, \dots, \beta_p$  are chosen using **the least square method**, which minimizes the sum of squared residuals (RSS):

$$RSS = \sum (y_i - \hat{y}_i)^2$$

where:

- $\sum$ : A greek symbol that means *sum*
- $y_i$ : The actual response value for the  $i$ th observation
- $\hat{y}_i$ : The predicted response value based on the multiple linear regression model

$$b_1 = \frac{(\sum x_2^2)(\sum x_1 y) - (\sum x_1 x_2)(\sum x_2 y)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)}$$

$$b_2 = \frac{(\sum x_1^2)(\sum x_2 y) - (\sum x_1 x_2)(\sum x_1 y)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)}$$

$$a = b_0 = \bar{Y} - b_1 \bar{X}_1 - b_2 \bar{X}_2$$

Example: Multiple Linear Regression by Hand

Suppose we have the following dataset with one response variable  $y$  and two predictor variables  $X_1$  and  $X_2$ :

$y$	$X_1$	$X_2$
140	60	22
155	62	25
159	67	24
179	70	20
192	71	15
200	72	14
212	75	14
215	78	11

Use the following steps to fit a multiple linear regression model to this dataset.

**Step 1: Calculate  $\sum X_1^2$ ,  $\sum X_2^2$ ,  $\sum X_1 y$ ,  $\sum X_2 y$  and  $\sum X_1 X_2$ .**

	y	X <sub>1</sub>	X <sub>2</sub>		X <sub>1</sub> <sup>2</sup>	X <sub>2</sub> <sup>2</sup>	X <sub>1</sub> y	X <sub>2</sub> y	X <sub>1</sub> X <sub>2</sub>
	140	60	22		3600	484	8400	3080	1320
	155	62	25		3844	625	9610	3875	1550
	159	67	24		4489	576	10653	3816	1608
	179	70	20		4900	400	12530	3580	1400
	192	71	15		5041	225	13632	2880	1065
	200	72	14		5184	196	14400	2800	1008
	212	75	14		5625	196	15900	2968	1050
	215	78	11		6084	121	16770	2365	858
Mean	181.5	69.375	18.125	Sum	38767	2823	101895	25364	9859
Sum	1452	555	145						

## Step 2: Calculate Regression Sums.

Next, make the following regression sum calculations:

- $\Sigma x_1^2 = \Sigma X_1^2 - (\Sigma X_1)^2 / n = 38,767 - (555)^2 / 8 = \mathbf{263.875}$
- $\Sigma x_2^2 = \Sigma X_2^2 - (\Sigma X_2)^2 / n = 2,823 - (145)^2 / 8 = \mathbf{194.875}$
- $\Sigma x_1y = \Sigma X_1y - (\Sigma X_1 \Sigma y) / n = 101,895 - (555 * 1,452) / 8 = \mathbf{1,162.5}$
- $\Sigma x_2y = \Sigma X_2y - (\Sigma X_2 \Sigma y) / n = 25,364 - (145 * 1,452) / 8 = \mathbf{-953.5}$
- $\Sigma x_1x_2 = \Sigma X_1X_2 - (\Sigma X_1 \Sigma X_2) / n = 9,859 - (555 * 145) / 8 = \mathbf{-200.375}$

	y	X <sub>1</sub>	X <sub>2</sub>		X <sub>1</sub> <sup>2</sup>	X <sub>2</sub> <sup>2</sup>	X <sub>1</sub> y	X <sub>2</sub> y	X <sub>1</sub> X <sub>2</sub>
	140	60	22		3600	484	8400	3080	1320
	155	62	25		3844	625	9610	3875	1550
	159	67	24		4489	576	10653	3816	1608
	179	70	20		4900	400	12530	3580	1400
	192	71	15		5041	225	13632	2880	1065
	200	72	14		5184	196	14400	2800	1008
	212	75	14		5625	196	15900	2968	1050
	215	78	11		6084	121	16770	2365	858
Mean	181.5	69.375	18.125	Sum	38767	2823	101895	25364	9859
Sum	1452	555	145						

Reg Sums	263.875	194.875	1162.5	-953.5	-200.375
----------	---------	---------	--------	--------	----------

## Step 3: Calculate b0, b1, and b2.

The formula to calculate  $b_1$  is:  $[(\sum x_2^2)(\sum x_1 y) - (\sum x_1 x_2)(\sum x_2 y)] / [(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2]$

Thus,  $b_1 = [(194.875)(1162.5) - (-200.375)(-953.5)] / [(263.875)(194.875) - (-200.375)^2] =$   
**3.148**

The formula to calculate  $b_2$  is:  $[(\sum x_1^2)(\sum x_2 y) - (\sum x_1 x_2)(\sum x_1 y)] / [(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2]$

Thus,  $b_2 = [(263.875)(-953.5) - (-200.375)(1152.5)] / [(263.875)(194.875) - (-200.375)^2] =$   
**-1.656**

The formula to calculate  $b_0$  is:  $y - b_1 X_1 - b_2 X_2$

Thus,  $b_0 = 181.5 - 3.148(69.375) - (-1.656)(18.125) =$  **-6.867**

**Step 5: Place  $b_0$ ,  $b_1$ , and  $b_2$  in the estimated linear regression equation.**

The estimated linear regression equation is:  $\hat{y} = b_0 + b_1 x_1 + b_2 x_2$

In our example, it is  $\hat{y} =$  **-6.867 + 3.148 $x_1$  - 1.656 $x_2$**

How to Interpret a Multiple Linear Regression Equation

Here is how to interpret this estimated linear regression equation:  $\hat{y} = -6.867 + 3.148x_1 - 1.656x_2$

**$b_0 = -6.867$ .** When both predictor variables are equal to zero, the mean value for  $y$  is -6.867.

**$b_1 = 3.148$ .** A one unit increase in  $x_1$  is associated with a 3.148 unit increase in  $y$ , on average, assuming  $x_2$  is held constant.

**$b_2 = -1.656$ .** A one unit increase in  $x_2$  is associated with a 1.656 unit decrease in  $y$ , on average, assuming  $x_1$  is held constant.