



## Module 3: Multivariate Linear Regression

Multivariate Regression is a supervised machine learning algorithm involving multiple data variables for analysis. Multivariate regression is an extension of multiple regression with one dependent variable and multiple independent variables. Based on the number of independent variables, we try to predict the output.

Multivariate regression tries to find out a formula that can explain how factors in variables respond simultaneously to changes in others.

There are numerous areas where multivariate regression can be used. Let's look at some examples to understand multivariate regression better.

1. Praneeta wants to estimate the price of a house. She will collect details such as the location of the house, number of bedrooms, size in square feet, amenities available, or not. Basis these details price of the house can be predicted and how each variables are interrelated.
2. An agriculture scientist wants to predict the total crop yield expected for the summer. He collected details of the expected amount of rainfall, fertilizers to be used, and soil conditions. By building a Multivariate regression model scientists can predict his crop yield. With the crop yield, the scientist also tries to understand the relationship among the variables.
3. If an organization wants to know how much it has to pay to a new hire, they will take into account many details such as education level, number of experience, job location, has niche skill or not. Basis this information salary of an employee can be predicted, how these variables help in estimating the salary.
4. Economists can use Multivariate regression to predict the GDP growth of a state or a country based on parameters like total amount spent by consumers, import expenditure, total gains from exports, total savings, etc.
5. A company wants to predict the electricity bill of an apartment, the details needed here are the number of flats, the number of appliances in usage, the number of people at home, etc. With the help of these variables, the electricity bill can be predicted.

The above example uses Multivariate regression, where we have many independent variables and a single dependent variable.



## Mathematical equation

The simple regression linear model represents a straight line meaning  $y$  is a function of  $x$ . When we have an extra dimension ( $z$ ), the straight line becomes a plane.

Here, the plane is the function that expresses  $y$  as a function of  $x$  and  $z$ . The linear regression equation can now be expressed as:

$$y = m_1.x + m_2.z + c$$

$y$  is the dependent variable, that is, the variable that needs to be predicted.  
 $x$  is the first independent variable. It is the first input.

$m_1$  is the slope of  $x$ . It lets us know the angle of the line ( $x$ ).

$z$  is the second independent variable. It is the second input.

$m_2$  is the slope of  $z$ . It helps us to know the angle of the line ( $z$ ).

$c$  is the intercept. A constant that finds the value of  $y$  when  $x$  and  $z$  are 0.

The equation for a model with two input variables can be written as:

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

What if there are three variables as inputs? Human visualizations can be only three dimensions. In the machine learning world, there can be  $n$  number of dimensions. The equation for a model with three input variables can be written as:

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

Below is the generalized equation for the multivariate regression model-

$$y = \beta_0 + \beta_1.x_1 + \beta_2.x_2 + \dots + \beta_n.x_n$$

Where  $n$  represents the number of independent variables,  $\beta_0 \sim \beta_n$  represents the coefficients, and  $x_1 \sim x_n$  is the independent variable.



The multivariate model helps us in understanding and comparing coefficients across the output. Here, the small cost function makes Multivariate linear regression a better model.

## Example: Multiple Linear Regression

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  $X_1^2$ ,  $X_2^2$ ,  $X_1y$ ,  $X_2y$  and  $X_1X_2$ .**

$y$	$X_1$	$X_2$
140	60	22
155	62	25
159	67	24
179	70	20
192	71	15
200	72	14
212	75	14
215	78	11
Mean	181.5	69.375
Sum	1452	555

$X_1^2$	$X_2^2$	$X_1y$	$X_2y$	$X_1X_2$
3600	484	8400	3080	1320
3844	625	9610	3875	1550
4489	576	10653	3816	1608
4900	400	12530	3580	1400
5041	225	13632	2880	1065
5184	196	14400	2800	1008
5625	196	15900	2968	1050
6084	121	16770	2365	858
Sum	38767	2823	101895	25364

**Step 2: Calculate Regression Sums.**

Next, make the following regression sum calculations:



Semester : VI

Subject : Machine Learning

Academic Year: 2023 - 2024

- $\Sigma X_1^2 = \Sigma X_1^2 - (\Sigma X_1)^2 / n = 38,767 - (555)^2 / 8 = \mathbf{263.875}$
- $\Sigma X_2^2 = \Sigma X_2^2 - (\Sigma X_2)^2 / n = 2,823 - (145)^2 / 8 = \mathbf{194.875}$
- $\Sigma X_1 Y = \Sigma X_1 Y - (\Sigma X_1 \Sigma Y) / n = 101,895 - (555 * 1,452) / 8 = \mathbf{1,162.5}$
- $\Sigma X_2 Y = \Sigma X_2 Y - (\Sigma X_2 \Sigma Y) / n = 25,364 - (145 * 1,452) / 8 = \mathbf{-953.5}$
- $\Sigma X_1 X_2 = \Sigma X_1 X_2 - (\Sigma X_1 \Sigma X_2) / n = 9,859 - (555 * 145) / 8 = \mathbf{-200.375}$

y	X <sub>1</sub>	X <sub>2</sub>
140	60	22
155	62	25
159	67	24
179	70	20
192	71	15
200	72	14
212	75	14
215	78	11
Mean	181.5	18.125
Sum	1452	145

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>
3600	484	8400	3080	1320
3844	625	9610	3875	1550
4489	576	10653	3816	1608
4900	400	12530	3580	1400
5041	225	13632	2880	1065
5184	196	14400	2800	1008
5625	196	15900	2968	1050
6084	121	16770	2365	858
Sum	38767	101895	25364	9859

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

### Step 3: Calculate b<sub>0</sub>, b<sub>1</sub>, and b<sub>2</sub>.

The formula to calculate b<sub>1</sub> is:  $[(\Sigma X_2^2)(\Sigma X_1 Y) - (\Sigma X_1 X_2)(\Sigma X_2 Y)] / [(\Sigma X_1^2)(\Sigma X_2^2) - (\Sigma X_1 X_2)^2]$

Thus, **b<sub>1</sub>** =  $[(194.875)(1162.5) - (-200.375)(-953.5)] / [(263.875)(194.875) - (-200.375)^2] = \mathbf{3.148}$

The formula to calculate b<sub>2</sub> is:  $[(\Sigma X_1^2)(\Sigma X_2 Y) - (\Sigma X_1 X_2)(\Sigma X_1 Y)] / [(\Sigma X_1^2)(\Sigma X_2^2) - (\Sigma X_1 X_2)^2]$

Thus, **b<sub>2</sub>** =  $[(263.875)(-953.5) - (-200.375)(1162.5)] / [(263.875)(194.875) - (-200.375)^2] = \mathbf{-1.656}$

The formula to calculate b<sub>0</sub> is:  $y - b_1 X_1 - b_2 X_2$

Thus, **b<sub>0</sub>** =  $181.5 - 3.148(69.375) - (-1.656)(18.125) = \mathbf{-6.867}$

### Step 5: Place b<sub>0</sub>, b<sub>1</sub>, and b<sub>2</sub> 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  $\hat{y} = \mathbf{-6.867 + 3.148x_1 - 1.656x_2}$

## How to Interpret a Multiple Linear Regression Equation

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



Semester : VI

Subject : Machine Learning

Academic Year: 2023 - 2024

**$b_0 = -6.867$** . 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.