

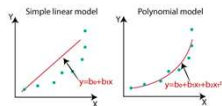
Multiple Linear regression

Assumptions

- **Independence:** the scores of any particular subject are independent of the scores of all other subjects
- **Normality:** in the population, the scores on the dependent variable are normally distributed for each of the possible combinations of the level of the X variables; each of the variables is normally distributed
- **Homoscedasticity:** in the population, the variances of the dependent variable for each of the possible combinations of the levels of the X variables are equal.
(In statistics, homoscedasticity is when a sequence of random variables has the same finite variance, or homogeneity of variance. It can also refer to when the variance in scores on one variable is similar at all values of another variable.)
- **Linearity:** In the population, the relation between the dependent variable and the independent variable is linear when all the other independent variables are held constant.

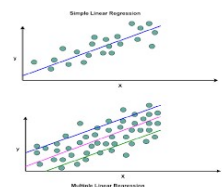
Simple vs. Multiple Regression

- One dependent variable Y predicted from one independent variable X
- One regression coefficient
- r^2 : proportion of variation in dependent variable Y predictable from X
- One dependent variable Y predicted from a set of independent variables (X_1, X_2, \dots, X_k)
- One regression coefficient for each independent variable
- R^2 : proportion of variation in dependent variable Y predictable by set of independent variables (X 's)



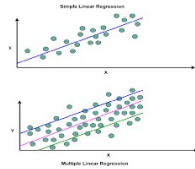
Multiple Regression Analysis (MRA)

- Method for studying the relationship between a dependent variable and two or more independent variables.
- Purposes:
 - Prediction
 - Explanation
 - Theory building



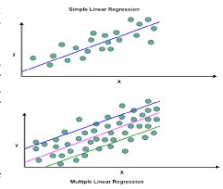
Operation?

- Uses the ordinary least squares solution (as does simple linear or bi-variable regression)
- Describes a line for which the (sum of squared) differences between the predicted and the actual values of the dependent variable are at a minimum.
- Represents the "function" that minimizes the sum of the squared errors.
- $Y_{\text{pred}} = a + a_1X_1 + a_2X_2 \dots + a_nX_n$



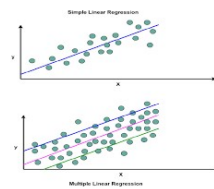
Operation?

- MLR produces a model that identifies the best weighted combination of independent variables to predict the dependent (or criterion) variable.
- $Y_{\text{pred}} = a + a_1X_1 + a_2X_2 \dots + a_nX_n$
 - MLR estimates the relative importance of several hypothesized predictors.
 - MLR assess the contribution of the combined variables to change the dependent variable.



Design Requirements

- One dependent variable (criterion)
- Two or more independent variables (predictor variables).
- Sample size: ≥ 50 (at least 10 times as many cases as independent variables)



Intuitive Method

For the one variable case, the calculation of b and a was:

$$a_1 = \frac{\sum xy}{\sum x^2}$$

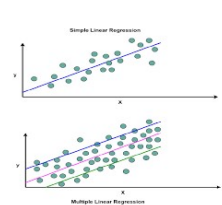
$$a = \bar{Y} - b\bar{X}$$

For the two variable case:

$$a_1 = \frac{(\sum x_1^2)(\sum x_2 y) - (\sum x_1 x_2)(\sum x_2 y)}{(\sum x_1^2)(\sum x_1^2) - (\sum x_1 x_2)^2}$$

and

$$a_2 = \frac{(\sum x_1^2)(\sum x_2 y) - (\sum x_1 x_2)(\sum x_1 y)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2}$$



Matrix Formulation

Example: Multiple Linear Regression
Find the Multiple Linear Regression model for the following given data points.

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

Suppose we have the following dataset with *one response variable y* and *two predictor variables X₁ and X₂*.

Example: Multiple Linear Regression
Find the Multiple Linear Regression model for the following given data points.
Suppose we have the following dataset with *one response variable y* and *two predictor variables X₁ and X₂*.

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

Step 1: Calculate X₁², X₂², X₁y, X₂y and X₁X₂.

	y	X ₁	X ₂		X ₁ ²	X ₂ ²	X ₁ y	X ₂ y	X ₁ X ₂
	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						

Example: Multiple Linear Regression
Find the Multiple Linear Regression model for the following given data points.
Suppose we have the following dataset with *one response variable y* and *two predictor variables X₁ and X₂*.

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

Step 2: Calculate Regression Sums.

Next, make the following regression sum calculations:

- $\sum X_1^2 = \sum X_1^2 - (\sum X_1)^2 / n = 38,767 - (555)^2 / 8 = 263.875$
- $\sum X_2^2 = \sum X_2^2 - (\sum X_2)^2 / n = 2,823 - (145)^2 / 8 = 194.875$
- $\sum X_1 y = \sum X_1 y - (\sum X_1)(\sum y) / n = 101,895 - (555)(1,452) / 8 = 1,162.5$
- $\sum X_2 y = \sum X_2 y - (\sum X_2)(\sum y) / n = 25,364 - (145)(1,452) / 8 = 1,162.5$
- $\sum X_1 X_2 = \sum X_1 X_2 - (\sum X_1)(\sum X_2) / n = 9,859 - (555)(145) / 8 = -200.375$

	y	X ₁	X ₂		X ₁ ²	X ₂ ²	X ₁ y	X ₂ y	X ₁ X ₂
	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

Example: Multiple Linear Regression

Find the Multiple Linear Regression model for the following given data points.

Suppose we have the following dataset with **one response variable y** and **two predictor variables X₁ and X₂**:

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

Step 3: Calculate a₀, a₁, and a₂.

The formula to calculate a₁ 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, **a₁ = 3.148**

The formula to calculate a₂ 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, **a₂ = -1.656**

The formula to calculate a₀ is: $y - a_1 X_1 - a_2 X_2$
Thus, **a₀ = -6.867**

Example: Multiple Linear Regression

Find the Multiple Linear Regression model for the following given data points.

Suppose we have the following dataset with **one response variable y** and **two predictor variables X₁ and X₂**:

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

Step 5: Place b₀, b₁, and b₂ in the estimated linear regression equation.

The estimated linear regression equation is: $\hat{y} = a_0 + a_1 X_1 + a_2 X_2$
In our example, it is **$\hat{y} = -6.867 + 3.148X_1 - 1.656X_2$**

How to Interpret a Multiple Linear Regression Equation

Here is how to interpret this estimated regression equation: $\hat{y} = -6.867 + 3.148X_1 - 1.656X_2$

a₀ = -6.867. When both predictor variables are equal to zero, the mean value for y is -6.867.

a₁ = 3.148. A one unit increase in x₁ is associated with a 3.148 unit increase in y, on average, assuming x₂ is held constant.

a₂ = -1.656. A one unit increase in x₂ is associated with a 1.656 unit decrease in y, on average, assuming x₁ is held constant.

<https://online.stat.psu.edu/stat462/node/133/>

Example 1: Pastry Sweetness Data- A designed experiment is done to assess how moisture content and sweetness of a pastry product affect a taster's rating of the product ([pastry.txt](#)). In a designed experiment, the eight possible combinations of four moisture levels and two sweetness levels are studied. Two pastries are prepared and rated for each of the eight combinations, so the total sample size is n = 16. The y-variable is the rating of the pastry. The two x-variables are moisture and sweetness. The values (and sample sizes) of the x-variables were designed so that the x-variables were not correlated.

Rating	Moisture	Sweetness
64	4	2
73	4	4
61	4	2
76	4	4
72	6	2
80	6	4
71	6	2
83	6	4
83	8	2
89	8	4
86	8	2
93	8	4
88	10	2
95	10	4
94	10	2
100	10	4

There is a linear relationship between rating and moisture and there is also a sweetness difference. The results given in the following output are for three different regressions - separate simple regressions for each x-variable and a multiple regression that incorporates both x-variables.

Regression Analysis: Rating versus Moisture Regression Equation
Rating = 50.77 + 4.425 Moisture

Regression Analysis: Rating versus Sweetness Regression Equation
Rating = 68.63 + 4.38 Sweetness

Regression Analysis: Rating versus Moisture, Sweetness
Regression Equation
Rating = 37.65 + 4.425 Moisture + 4.375 Sweetness

Example 2: Female Stat Students

The data are from n = 214 females in statistics classes at the University of California at Davis ([stat_females.txt](#)). The variables are y = student's self-reported height, x1 = student's guess at her mother's height, and x2 = student's guess at her father's height. All heights are in inches. The scatterplots below are of each student's height versus mother's height and student's height against father's height.

Scatterplot of Height vs momheight, dadheight

Regression Equation
Height = 18.55 + 0.3035 momheight + 0.3879 dadheight

Example 2: Female Stat Students

The data are from $n = 214$ females in statistics classes at the University of California at Davis (stat_females.txt). The variables are y = student's self-reported height, x_1 = student's guess at her mother's height, and x_2 = student's guess at her father's height. All heights are in inches. The scatterplots below are of each student's height versus mother's height and student's height against father's height.

Interpretations

• The sample multiple regression equation is predicted student height = $18.55 + 0.3035 \times \text{mother's height} + 0.3879 \times \text{father's height}$. To use this equation for prediction, we substitute specified values for the two parents' heights.

• We can interpret the "slopes" in the same way that we do for a simple linear regression model but we have to add the constraint that values of other variables remain constant. For example:

- When father's height is held constant, the average student height increases 0.3035 inches for each one-inch increase in mother's height.
- When mother's height is held constant, the average student height increases 0.3879 inches for each one-inch increase in father's height.