

# What is Multiple Linear Regression?

Authored by  
**stats writer**

December 19, 2025

## RECOMMENDED CITATION

stats writer (2025). *What is Multiple Linear Regression?*. PSYCHOLOGICAL SCALES.  
Retrieved from <https://scales.arabpsychology.com/?p=107956>

Multiple linear regression (MLR) is a powerful statistical technique utilized across various scientific and business disciplines to model the linear relationship between several explanatory factors and a single outcome variable. Unlike simple linear regression, which relies on only one predictor, MLR allows researchers and analysts to account for the simultaneous influence of two or more independent variables on a specific dependent variable. This comprehensive approach provides a much richer understanding of complex systems.

The fundamental goal of an MLR model is twofold: first, to determine the strength and direction of the association between each predictor and the outcome, and second, to generate a predictive equation that can reliably estimate the value of the outcome variable given a set of known predictor values. This capability makes MLR indispensable for forecasting and for performing crucial control functions, allowing practitioners to understand how changes in various predictors might collectively affect the response.

## The Distinction Between Simple and Multiple Linear Regression

When the research focus is narrow--examining the association between a single predictor variable and a response variable--the appropriate tool is often simple linear regression. This foundational method provides initial insights into bivariate relationships, offering a baseline understanding of how two factors interact.

However, real-world phenomena are rarely governed by a single factor. To capture the complexity of outcomes influenced by several factors simultaneously, we pivot to multiple linear regression. MLR allows us to isolate the unique contribution of each predictor while controlling for the effects of others, leading to much more robust and realistic models necessary for empirical research and data science applications.

## The Core Mathematical Model

If we incorporate  $p$  distinct predictor variables into our analysis, the mathematical expression for the multiple linear regression model adheres to a standard linear form. This equation represents the theoretical relationship between the expected value of the response and the observed values of the predictors:

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

Understanding each component of this equation is critical for proper model interpretation and application. The definitions below clarify the role of each term in the linear framework:

**Y:** This represents the response variable, which is the outcome we are attempting to predict or explain.

**X<sub>j</sub>**: This denotes the jth predictor variable (or independent variable) used in the model.

**β<sub>0</sub>**: Known as the Y-intercept, this is the expected value of Y when all predictor variables (X<sub>1</sub> through X<sub>p</sub>) are equal to zero.

**β<sub>j</sub>**: This is the regression coefficient associated with the jth predictor X<sub>j</sub>. It quantifies the estimated average change in Y resulting from a one-unit increase in X<sub>j</sub>, crucially, while holding the values of all other predictors fixed.

**ε**: This is the error term (or residual term). It accounts for the unexplained variation in Y that the model cannot capture, representing the difference between the actual observed value and the value predicted by the linear relationship.

## Estimating Coefficients: The Least Squares Method

The process of fitting the MLR model to a dataset involves estimating the unknown population parameters (the  $\beta$  coefficients). These coefficient estimates (often denoted as  $\hat{\beta}$ ) are calculated using a process known as the least squares method (OLS). The core principle of OLS is to find the line (or hyperplane, in the case of MLR) that best fits the data by minimizing the overall prediction error, thereby generating the line of best fit.

Specifically, the least squares method minimizes the Sum of Squared Residuals (RSS). The RSS represents the total discrepancy between the observed data points and the values predicted by the model. The mathematical formula defining the quantity that OLS seeks to minimize is:

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

Within this formulation, each element plays a precise role in quantifying the error:

**Σ**: This Greek symbol represents the operation of summation, summing the squared errors across all data points.

**y<sub>i</sub>**: This is the actual observed response value for the ith data point.

**ŷ<sub>i</sub>**: This is the predicted response value for the ith observation, calculated using the estimated multiple linear regression equation.

While the underlying calculation of these coefficients relies on intricate matrix algebra--a topic typically reserved for advanced statistical courses--modern statistical software handles these computational details automatically. Programs like R, Python (with libraries like scikit-learn or Statsmodels), or specialized statistical packages calculate the optimal coefficients instantaneously, allowing researchers to focus on model specification and interpretation rather than manual calculation.

## Interpreting Multiple Linear Regression Output

Once the MLR model is computed, the next crucial step is interpreting the coefficients and statistical metrics provided in the output summary. Let us consider a classic scenario: modeling a student's final exam score (the response variable) based on two predictors: the total number of *hours studied* and the number of *prep exams taken*.

The output generated by statistical software packages (such as R, SPSS, or Excel, as shown below) summarizes the results, providing the estimated coefficients, standard errors, t-statistics, and associated p-values for each variable. This comprehensive summary forms the basis for constructing the predictive equation and assessing the significance of individual predictors.

**Note:** The provided screenshot illustrates a typical output structure. While this specific image refers to multiple linear regression output for Excel, the core metrics and numerical values shown are representative of standard regression analysis across any major statistical platform.

D	E	F	G	H	I	J	K
SUMMARY OUTPUT							
<i>Regression Statistics</i>							
Multiple R	0.857						
R Square	0.734						
Adjusted R Square	0.703						
Standard Error	5.366						
Observations	20						
ANOVA							
	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>Significance F</i>		
Regression	2	1350.76	675.38	23.46	0.00		
Residual	17	489.44	28.79				
Total	19	1840.20					
	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>	
Intercept	67.67	2.82	24.03	0.00	61.73	73.61	
hours	5.56	0.90	6.18	0.00	3.66	7.45	
prep_exams	-0.60	0.91	-0.66	0.52	-2.53	1.33	

Based on the coefficients derived from this model output, we can formulate the estimated multiple linear regression equation, which mathematically relates the predictors to the predicted exam score.

## Understanding the Regression Coefficients

The interpretation of MLR coefficients (the  $\beta$  values) is contingent upon the assumption of "all else being equal" (or **ceteris paribus**). This means that the effect of one variable is assessed independently of the others, assuming their values are held constant throughout the change.

For our example, the interpretation of the estimated coefficients (67.67 for the intercept, 5.56 for hours studied, and -0.60 for prep exams taken) is detailed as follows:

**Hours Studied Coefficient (5.56):** Each additional unit increase in hours studied is associated with an average increase of **5.56** points in the final exam score. This interpretation is valid only if the number of *prep exams taken* is assumed to be held constant across all students being compared.

**Prep Exams Taken Coefficient (-0.60):** Conversely, each additional prep exam taken is associated with an average decrease of **0.60** points in the final exam score. This relationship is observed assuming the total *hours studied* remains unchanged. The negative sign suggests an inverse relationship, which might warrant further investigation into the data collection process or underlying variables.

**Intercept (67.67):** This is the predicted exam score for a student who has studied 0 hours and taken 0 prep exams.

This estimated model provides a practical tool for prediction. We can input specific predictor values to forecast an expected outcome. For instance, if a student studies for 4 hours ( $X_1=4$ ) and takes 1 prep exam ( $X_2=1$ ), the expected exam score is calculated as:

$$\text{Exam score} = 67.67 + 5.56*(4) - 0.60*(1) = \mathbf{89.31}$$

## Evaluating Model Significance and Fit

Beyond the individual coefficients, the regression output provides crucial statistics that evaluate the overall effectiveness and reliability of the model. Assessing the model's significance helps determine if the set of predictors, taken together, contribute meaningfully to explaining the variation in the response variable.

**F Statistic and Significance F:** The overall F statistic tests the null hypothesis that all regression coefficients are zero. The associated global p-value (Significance F) indicates whether the model as a whole is statistically significant. In our example, a p-value less than 0.05 suggests that the variables *hours studied* and *prep exams taken* collectively have a significant association with the *exam score*.

**Coefficient P-values:** These values assess the significance of each predictor individually. We see that *hours studied* ( $p = 0.00$ ) is highly significant, while *prep exams taken* ( $p = 0.52$ ) is not

statistically significant at  $\alpha = 0.05$ . This finding often leads analysts to consider removing the non-significant variable to simplify the model.

### R-Squared (Coefficient of Determination)

The R-Squared value, also known as the coefficient of determination, provides a relative measure of model fit. Specifically, it represents the proportion of the total variance in the response variable that is explained by the full set of predictor variables included in the model. Its value is always bounded between 0 and 1. In this example, 73.4% of the variation in the exam scores can be collectively attributed to the variance in hours studied and prep exams taken.

### Standard Error of the Estimate

The Standard Error of the Estimate provides an absolute measure of model error, measured in the same units as the dependent variable. It represents the average distance that the observed data points fall from the estimated regression line. The smaller the standard error, the better the predictive capability. In our student example, a Standard Error of 5.366 units means that, on average, the model's prediction for an exam score will be off by approximately 5.366 points. If making predictions is the primary goal, this metric is often more useful than R-Squared.

For a more detailed discussion regarding the comparative advantages of R-Squared versus the Standard Error when evaluating model fit and predictive capability, consult the following resources:

[What is a Good R-squared Value?](#)

[Understanding the Standard Error of a Regression Model](#)

### Key Assumptions of Multiple Linear Regression

For the results of a multiple linear regression analysis to be valid, reliable, and unbiased, the underlying data structure and the calculated residuals must satisfy several critical assumptions. Violating these assumptions can lead to incorrect coefficient estimates, unreliable p-values, and misleading standard errors.

**Linearity:** The relationship between the independent variables and the mean of the dependent variable must be linear. This means the model equation should accurately reflect the underlying structure without requiring complex non-linear transformations.

**Independence of Errors:** The residuals (or errors) must be statistically independent of one another. This is crucial, especially for time-series data, where autocorrelation between consecutive errors must be absent.

**Homoscedasticity:** The variance of the residuals must remain constant across all levels of the predictor variables. If the variability of the errors increases or decreases systematically, the

condition of heteroscedasticity is met, which biases the standard errors.

**Normality of Errors:** The residuals of the model should be approximately normally distributed. While this assumption is less strict for large datasets, severe non-normality can impact the accuracy of inferential statistics.

**No Multicollinearity:** Although not always listed as a core statistical assumption of the error term, high correlation among the predictor variables (multicollinearity) can destabilize the coefficient estimates, making them difficult to interpret and highly sensitive to small changes in the input data.

A thorough understanding of diagnostic plots and formal tests is required to confirm that these conditions are met before drawing definitive conclusions from the regression output. For a comprehensive guide on how to test and address potential violations of these prerequisites, refer to [this detailed article](#).

## Implementing Multiple Linear Regression Using Statistical Tools

While the theoretical background of MLR involves complex mathematics, the execution and analysis are simplified greatly by modern statistical software. These tools efficiently handle the extensive matrix calculations required by the least squares method and provide comprehensive output summaries necessary for interpretation and diagnostics.

A wide range of platforms, from dedicated statistical packages to programming environments, support MLR analysis. The tutorials below offer detailed, step-by-step guidance on implementing the model across various popular software solutions:

[How to Perform Multiple Linear Regression in R](#)

[How to Perform Multiple Linear Regression in Python](#)

[How to Perform Multiple Linear Regression in Excel](#)

[How to Perform Multiple Linear Regression in SPSS](#)

[How to Perform Multiple Linear Regression in Stata](#)

[How to Perform Linear Regression in Google Sheets](#)