

How to Perform Simple Linear Regression and Understand Its Results

Authored by
stats writer

January 22, 2026

RECOMMENDED CITATION

stats writer (2026). *How to Perform Simple Linear Regression and Understand Its Results*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=127124>

Simple linear regression (SLR) is a fundamental statistical method used extensively across diverse fields such as economics, psychology, finance, and engineering. Its primary function is to model and analyze the linear relationship between two continuous variables: a single dependent variable (the outcome being predicted) and one independent variable (the predictor). This powerful technique involves fitting the best possible straight line through a scatter plot of observed data points, allowing researchers to quantify the strength and determine the direction of the relationship. By establishing this mathematical connection, Simple Linear Regression provides critical insights, enabling accurate predictions of the dependent variable's value based on a given change in the independent variable. Understanding SLR is essential for anyone engaged in data analysis, as it forms the bedrock for more complex multivariate statistical models.

What is Simple Linear Regression?

Simple Linear Regression is a statistical modeling technique dedicated to forecasting the value of an outcome variable using just one predictor variable. It achieves this by calculating the equation for the line that minimizes the sum of squared errors between the line and the actual data points--this is often referred to as the method of least squares. The resulting regression equation, often expressed as $Y = \beta_0 + \beta_1 X + \epsilon$, provides a precise numerical quantification of the relationship between the two variables, where Y is the dependent variable, X is the independent variable, β_0 is the intercept, β_1 is the slope, and ϵ represents the error term.

A fundamental requirement for applying Simple Linear Regression successfully is that the variable you intend to predict (the outcome) must be continuous. Continuous variables are those that can take on any value within a given range, such as temperature, height, or monetary value. Furthermore, the data utilized in the analysis must rigorously meet a specific set of underlying statistical assumptions. Failing to satisfy these critical assumptions can lead to unreliable models, biased estimates, and potentially erroneous conclusions about the underlying population relationship.

Simple Linear Regression



Assumptions for Simple Linear Regression

Every statistical modeling approach is built upon foundational assumptions. These assumptions dictate specific properties that your dataset must exhibit in order for the statistical results derived from the method to be statistically valid, unbiased, and generalizable. When applying Simple Linear Regression, careful verification of these assumptions is mandatory prior to interpreting the model's coefficients or relying on its predictive capabilities. Violation of these assumptions necessitates either data transformation or the use of alternative, more robust statistical techniques.

The primary assumptions that must be satisfied for a robust Simple Linear Regression model are critical for ensuring the fidelity and accuracy of the analysis. Researchers must meticulously test for each of these conditions, often using diagnostic plots and specific statistical tests, to confirm that the selected model is appropriate for the data structure being analyzed. These assumptions directly influence the validity of the hypothesis tests (like p-values) and the confidence intervals derived from the model.

The core assumptions for Simple Linear Regression include:

Linearity of the Relationship

Absence of Significant Outliers
Homoscedasticity (Similar Spread across Range)
Independence of Observations
Normality of the Residuals

Let's delve into the specifics of each of these crucial requirements to fully understand their implications for model reliability and data structure.

Linearity of the Relationship

The assumption of linearity requires that the relationship between the independent variable and the dependent variable is adequately represented by a straight line. If the true relationship is curvilinear (e.g., quadratic or exponential), fitting a linear model will introduce systematic error and lead to poor predictions. To assess linearity, researchers typically examine a scatter plot of the two variables. If the data points visually cluster around a hypothetical straight line, the assumption is generally considered met.

It is important to note that this assumption does not imply that the data points must fall perfectly on a line, but rather that the underlying trend is linear. If non-linearity is detected, transformation techniques (such as logarithmic or square root transformations) may be applied to one or both variables to linearize the relationship before proceeding with the regression analysis. Alternatively, more complex regression models, such as polynomial regression, may be required.

Absence of Significant Outliers

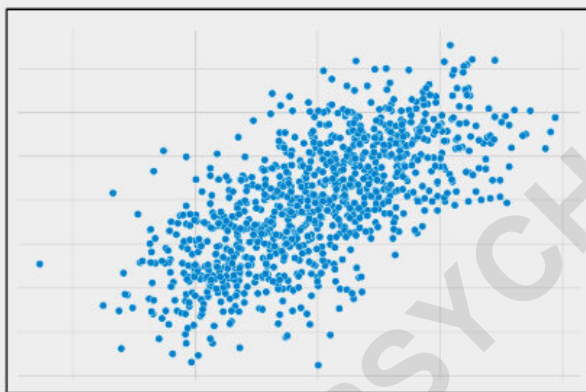
Linear Regression models are highly sensitive to outliers, which are data points possessing unusually large or small values relative to the rest of the dataset. An outlier, especially one that is highly influential, can drastically skew the slope (β) and intercept (β) of the fitted line, leading to a model that poorly represents the majority of the data. This is because the least squares method works by minimizing squared errors, meaning large errors caused by outliers are disproportionately weighted.

Detection of outliers is typically achieved through visual inspection of scatter plots or residual plots, as well as by calculating diagnostic statistics such as Cook's Distance or leverage scores. When outliers are identified, careful consideration must be given to their origin. If the outlier is due to a measurement error or data entry mistake, it should be corrected or removed. If it represents a genuinely extreme but valid observation, researchers might consider using robust regression techniques that are less sensitive to extreme values, rather than simply discarding the data point.

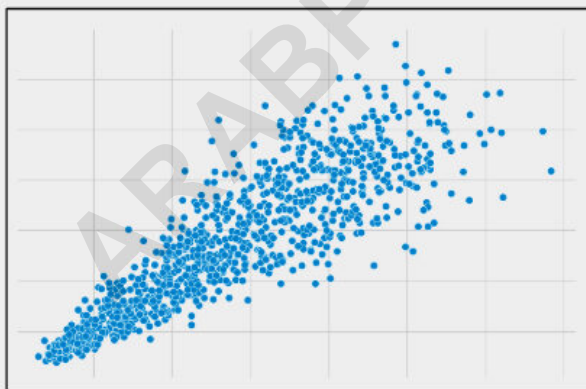
Homoscedasticity (Similar Spread across Range)

The term homoscedasticity refers to the consistency of the variance of the errors (residuals) across all levels of the independent variable. In simpler terms, it means the spread or variation of the data points around the regression line should remain relatively constant as the independent variable increases or decreases. If the spread of the residuals systematically changes (e.g., the spread widens as X increases), this condition is violated, resulting in a state called heteroscedasticity.

When heteroscedasticity is present, the standard errors of the regression coefficients become unreliable, leading to inaccurate hypothesis testing and confidence intervals. This violation does not bias the coefficient estimates themselves, but it does render statistical inference questionable. Visual inspection of a plot of the residuals versus the predicted values is the primary diagnostic tool. If a funnel or cone shape is observed, heteroscedasticity is likely present. Solutions include transforming the dependent variable or using weighted least squares regression methods.



These data have a similar spread across their range.



These data have greater spread at higher values.

Independence of Observations

The assumption of independence mandates that each observation (or data point) in the dataset must be independent of all other observations. This means that the value of the error term for one

data point should not be related to the error term of any other data point. This assumption is crucial because statistical inference relies on the idea that the errors are randomly and independently distributed. Violations of independence often occur in time-series data or clustered data structures.

A common scenario where this assumption is violated involves collecting multiple data points over time from the same unit of observation (e.g., the same subject, customer, or geographical location). Data collected repeatedly from the same source are inherently related (or correlated); for instance, a subject's blood pressure reading today is likely related to their reading yesterday. Such dependence invalidates standard linear regression assumptions. In these cases, researchers should avoid Simple Linear Regression and instead utilize more advanced longitudinal modeling techniques, such as a **Mixed Effects Model**, which explicitly accounts for the hierarchical or repeated nature of the data.

If your data have repeated measures over time from the same units of observation, you should use a Mixed Effects Model, which is specifically designed to handle dependent data structures.

Normality of Residuals

The final assumption focuses on the distribution of the residuals. Residuals are defined as the differences between the actual observed values of the dependent variable and the values predicted by the regression line. The assumption requires that the distribution of these error terms must follow a normal distribution, often visualized as a bell curve.

Meeting this assumption is particularly important for the accurate calculation of confidence intervals and p-values, especially in smaller sample sizes. When the residuals are normally distributed, it assures that the statistical inferences (like significance testing) are equally applicable across the full spread of the data and that there is no systematic bias in the prediction process. Researchers commonly assess this assumption using Q-Q plots, histograms of residuals, and formal statistical tests such as the Shapiro-Wilk test.

Choosing the Appropriate Model: When to use Simple Linear Regression?

Selecting the correct statistical model is the most critical step in data analysis. Simple Linear Regression is the appropriate choice only when the research question and the structure of your data align with three fundamental criteria. Using SLR outside of these criteria will yield inaccurate or misleading results, forcing the analyst to consider alternatives like logistic regression, time series analysis, or multivariate methods.

You should employ Simple Linear Regression only under the following specific conditions:

The goal is either to perform a **prediction** of an outcome or to precisely quantify the numerical linear relationship between two variables.

The variable you are attempting to predict (the dependent variable) is strictly **continuous**.

You have exactly **one independent variable**, or one single predictor, measured at one point in time.

Understanding the nuance of these criteria is essential for differentiating SLR from other common modeling techniques. Let's clarify these points to provide guidance on appropriate model selection.

The Goal is Prediction or Quantification

The primary utility of Simple Linear Regression is twofold: prediction and quantification. When you seek a statistical test capable of forecasting the value of one variable based on the input of another, you are asking a prediction question, which SLR is designed to answer. For instance, predicting future sales revenue based on current advertising expenditure falls perfectly within the scope of SLR, provided the other assumptions hold true. Beyond prediction, SLR provides a clear mathematical coefficient that quantifies the exact marginal change in the dependent variable resulting from a unit change in the independent variable.

It is crucial to distinguish this objective from other analytical goals. For example, if your sole purpose is examining the strength and direction of the linear association between two variables without implying causality or prediction, a simple correlation analysis (like Pearson's r) might suffice. If the goal is to examine differences in an outcome between predefined groups, an Analysis of Variance (ANOVA) or T-test would be more appropriate.

The Dependent Variable Must Be Continuous

As mentioned earlier, the variable being predicted must be continuous. Continuous variables are measured on a scale where theoretically infinite values are possible between any two points. Classic examples include physical measurements such as height, weight, time, or variables like standardized test scores or calculated ratios. These variables offer detailed numeric variation necessary for fitting a precise regression line.

It is vital to recognize data types that violate this requirement. These include ordered data (e.g., ranking systems, Likert scales treated as ordinal), categorical data (e.g., gender, eye color, type of treatment), or binary/dichotomous data (e.g., purchased the product or not, presence or absence of a disease). Applying SLR to these non-continuous outcomes will violate the normality and homoscedasticity assumptions, rendering the resulting model invalid and the interpretations meaningless.

*If your dependent variable is binary (two outcomes), you should use **Simple Logistic Regression**.*

If your dependent variable is categorical with three or more distinct classes, then you should consider **Multinomial Logistic Regression** or **Linear Discriminant Analysis**.

Requirement of One Independent Variable

The "Simple" in Simple Linear Regression strictly implies the use of only one predictor variable (the independent variable) to model the outcome. This predictor must also be measured cross-sectionally--that is, at a single, consistent point in time for all observations. SLR allows for an easy-to-interpret visualization of the relationship, as the model is merely a two-dimensional straight line (Y vs. X).

If your research design involves incorporating multiple predictor variables simultaneously to enhance the accuracy of your forecast (e.g., using both advertising dollars and competitor pricing to predict revenue), Simple Linear Regression is insufficient. Likewise, if your single independent variable is measured repeatedly for the same subjects over time, violating the independence assumption, you must select an alternative modeling approach.

If you have more than one independent variable acting as predictors, you should transition to a generalization of this method called **Multiple Linear Regression**. Furthermore, if you have one independent variable but it is measured for the same group at multiple points in time, you must use a **Mixed Effects Model** or a time series model to account for the resulting data dependence.

Simple Linear Regression Example and Interpretation

To illustrate the practical application of Simple Linear Regression, consider a common business scenario focused on optimizing marketing spend. We aim to determine if and how the amount spent on advertising influences sales revenue across different metropolitan areas.

In this study, the variables are defined as:

Dependent Variable: Revenue (measured in currency, which is continuous)

Independent Variable: Dollars spent on advertising per city (also measured in currency, and continuous)

Before running the model, we establish the null hypothesis (H_0), which represents the scenario if the predictor has no effect. The statistical null hypothesis is that there is no linear relationship between dollars spent on advertising and the resulting revenue within a city, meaning the slope of the line (β) is zero. Our regression analysis will assess the probability of observing our data if this null hypothesis were truly correct.

After meticulously gathering the data and confirming that all underlying statistical assumptions of

linear regression are met—including linearity, homoscedasticity, and normality of residuals—we execute the formal analysis. The output of this procedure yields key statistical metrics that allow for interpretation and decision-making.

Interpreting the Regression Coefficients (Beta and P-value)

The core results of a linear regression model are the beta coefficients (β). For simple linear regression, there are two coefficients: the intercept (β_0) and the slope (β_1). The intercept (β_0) represents the predicted value of the dependent variable when the independent variable is exactly zero. In our example, it would be the predicted revenue if zero dollars were spent on advertising.

The slope coefficient (β_1), also known as the regression weight, is the key parameter representing the numerical relationship between the variables. This coefficient indicates the expected change in the dependent variable (Revenue) for every one-unit increase in the independent variable (Advertising Dollars). For example, if $\beta_1 = 1.5$, this means that for every additional dollar spent on advertising, the revenue is expected to increase by \$1.50, holding all other factors constant.

Crucially, each coefficient is associated with a p-value. The p-value for β_1 tests the statistical significance of the relationship. It represents the probability of observing a slope as extreme or more extreme than β_1 if the true relationship in the population were zero (i.e., if the null hypothesis were true). A p-value less than or equal to the significance level (typically 0.05) indicates that the result is statistically significant, allowing us to reject the null hypothesis and confidently conclude that the relationship between advertising spend and revenue is unlikely to be due to random chance alone.

Understanding Model Fit (R-Squared)

In addition to the coefficients, the analysis generates an R-Squared (R^2) value, formally known as the Coefficient of Determination. This metric provides an assessment of the overall model fit. The R^2 value ranges from 0 to 1, and it quantifies the proportion of the variance in the dependent variable that is predictable from the independent variable.

A higher R^2 indicates a better fit; for example, an R^2 of 0.75 means that 75% of the variability in Revenue is explained by the variability in Advertising Dollars. Conversely, an R^2 close to zero suggests that the linear model provides little explanation for the outcome. While a high R^2 is desirable, researchers must interpret it cautiously, as a good fit does not guarantee that the causal assumptions are met, nor does it confirm the absence of bias from violations of other critical assumptions like independence or homoscedasticity.