

Understanding Residuals in Regression Analysis: A Step-by-Step Guide

Authored by
Mohammed Iooti

November 9, 2025

RECOMMENDED CITATION

Mohammed Iooti (2025). *Understanding Residuals in Regression Analysis: A Step-by-Step Guide*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=14022>

[Simple linear regression](#) is a foundational statistical method widely employed across scientific, economic, and business domains. Its fundamental goal is to mathematically model and quantify the relationship between two continuous variables: an independent factor, commonly represented as x , and a dependent outcome, designated as y . By successfully establishing this linear relationship, we gain the powerful ability to predict the value of the outcome variable based on the input of the independent variable.

In the context of regression analysis, the input variable x is formally known as the [predictor variable](#), as it is utilized to forecast changes in the other variable. Conversely, the outcome variable y , which represents the measured response we are attempting to model, is referred to as the [response variable](#). Constructing a successful regression model provides valuable, quantifiable insight into both the strength and the direction of the linear association between these two critical components.

Modeling Relationships Using Simple Linear Regression

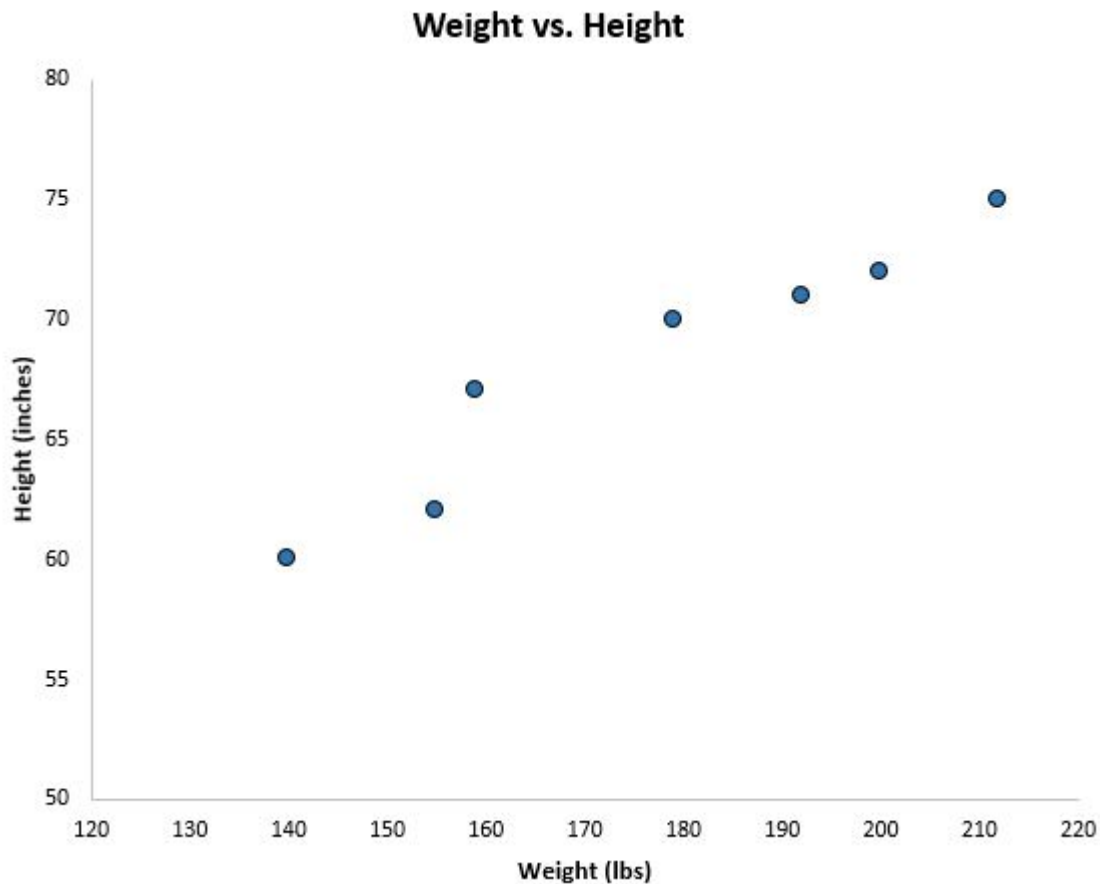
To practically illustrate the implementation of regression analysis, let us consider a small dataset containing anthropometric measurements. Suppose we have collected data detailing the weight (our potential [predictor variable](#), x) and the corresponding height (the [response variable](#), y) for a group of seven individuals. Analyzing this type of paired data allows us to rigorously determine if, and how strongly, weight correlates with or influences height within this specific sample population.

The raw data collected for this illustrative example is systematically presented below. For our analysis, we designate *weight* as the factor used for making predictions and *height* as the outcome we are attempting to model using a linear equation.

Weight (lbs)	Height (inches)
140	60
155	62
159	67
179	70
192	71
200	72
212	75

Prior to engaging in formal calculations, the standard practice in data analysis dictates that we first visualize the relationship using a [scatter plot](#). By mapping weight onto the horizontal (x) axis and height onto the vertical (y) axis, we can visually inspect the relationship between the variables. The resulting plot immediately reveals the overall trend, allowing us to confirm whether a linear model is

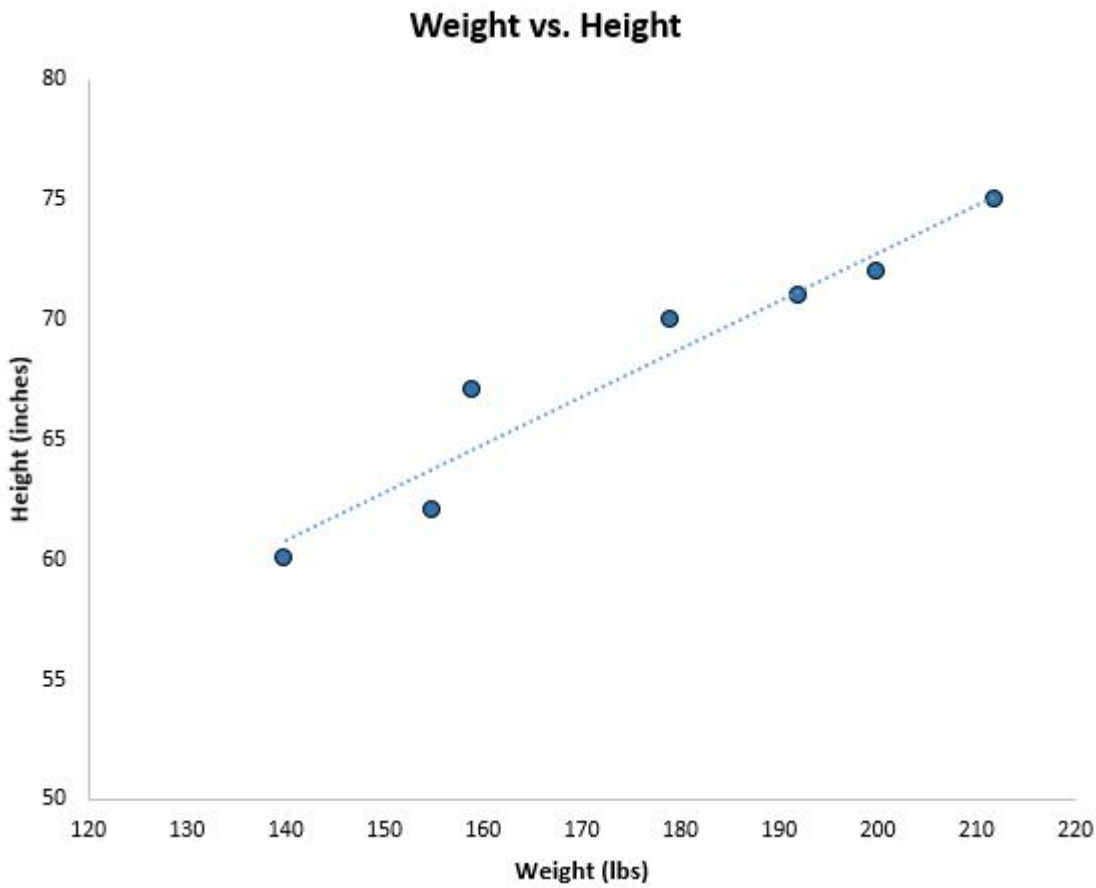
indeed the most appropriate structure for describing the data.



Calculating the Optimal Line of Best Fit

Visual inspection of the [scatter plot](#) above strongly suggests a positive correlation: as weight increases, height generally tends to increase as well. However, this preliminary visual confirmation is purely qualitative. To accurately **quantify** the precise strength and nature of the linear relationship between weight and height, we must apply the principles of [linear regression](#). This statistical procedure mathematically determines the single optimal straight line that best summarizes the observed trend within the entire dataset.

This calculated trend line is universally known as the [line of best fit](#). It is derived using a methodology called Ordinary Least Squares (OLS), a technique designed specifically to minimize the sum of the squared vertical distances between the line and every individual data point. This crucial minimization process yields the single most representative linear model for the dataset, which is visually depicted below cutting through the cloud of data points.



The mathematical formula that defines this [line of best fit](#) in [simple linear regression](#) is expressed as:

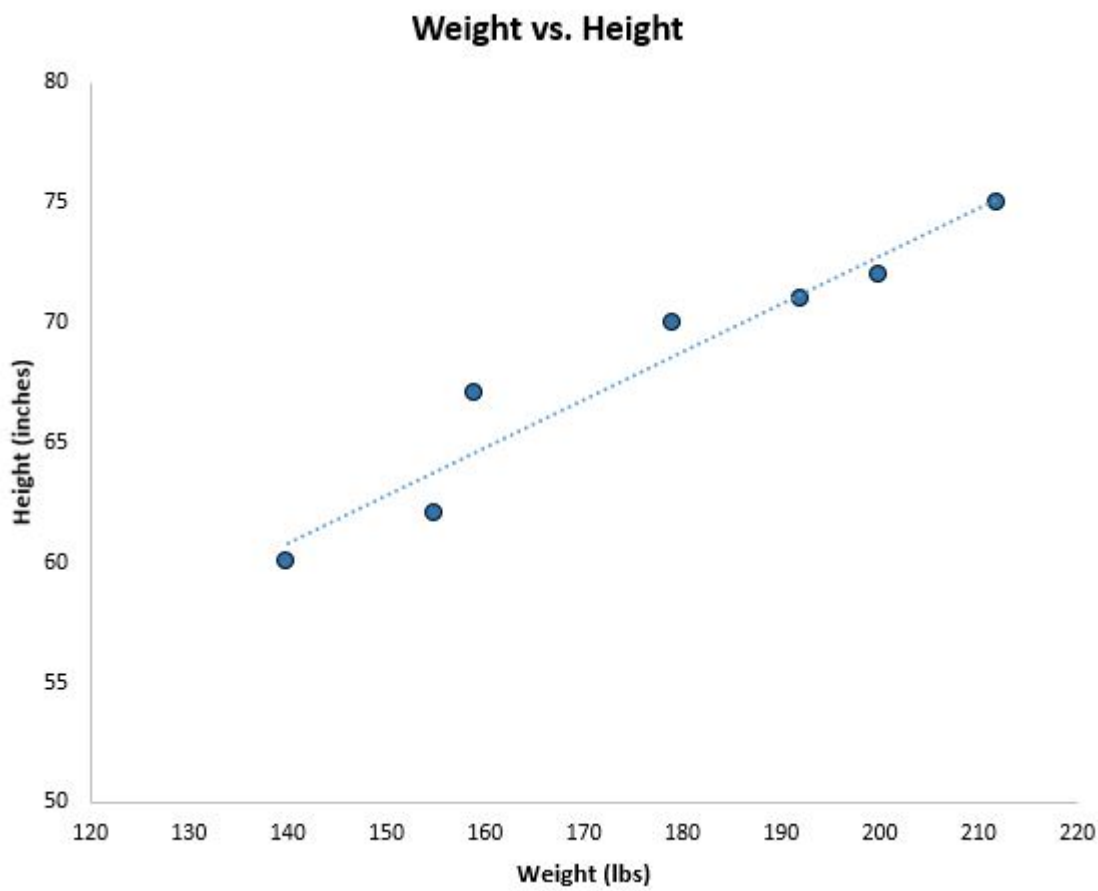
$$? = b_0 + b_1x$$

In this equation, ? represents the **predicted value** of the response variable (y); b_0 is the y-intercept (the predicted value of y when x is zero); b_1 is the [regression coefficient](#) (which quantifies the change in ? for every one-unit change in x); and x is the actual value of the [predictor variable](#). For our specific weight and height example, the necessary calculations resulted in the following derived model equation:

$$\text{height} = 32.783 + 0.2001 * (\text{weight})$$

Defining and Calculating Residuals: The Error Term

While the [line of best fit](#) provides the single optimal prediction for the overall population trend, a careful examination of the scatter plot confirms that individual data points rarely fall exactly upon this line:



This vertical distance between an observed data point (the actual y-value) and the value predicted by the regression line (?) is formally termed the **residual** (often denoted as 'e'). A residual is, fundamentally, the error in prediction specific to that single observation. For every individual data point in our dataset, we can calculate its corresponding residual by determining the difference between the actual observed response value (y) and the predicted response value (?) derived from the **line of best fit**.

The formal equation used for calculating a single **residual** is straightforward:

$$\text{Residual (e)} = \text{Actual Value (y)} - \text{Predicted Value (?)}$$

Understanding the sign of the residual is vital for interpreting the model's performance. A positive **residual** indicates that the actual observation lies **above** the regression line, meaning the model underestimated the true value. Conversely, a negative residual signifies that the observation lies **below** the line, which means the model overestimated the true value. In either case, the magnitude (absolute value) of the residual directly reflects the degree of error in the prediction for that specific data point.

Step-by-Step Residual Calculation Examples

We will now walk through the process of calculating the [residual](#) for two specific individuals using our derived regression equation: $\text{height} = 32.783 + 0.2001 * (\text{weight})$.

Example 1: Calculating the Residual for Individual 1

Referring back to our initial dataset, the measurements recorded for the first individual are a weight (x) of **140** lbs. and an actual height (y) of **60** inches.

Weight (lbs)	Height (inches)
140	60
155	62
159	67
179	70
192	71
200	72
212	75

To determine the predicted height (?) for this person, we substitute their weight ($x = 140$) into the established line of best fit equation:

$$\text{height} = 32.783 + 0.2001 * (\text{weight})$$

The calculation for the predicted height proceeds as follows:

$$\text{height} = 32.783 + 0.2001 * (140)$$

$$\text{height (?) = 60.797 inches}$$

We then compare the predicted height (60.797 inches) to the actual observed height (60 inches) to find the residual:

$$\text{Residual} = \text{Actual Height (y)} - \text{Predicted Height (?)}$$

The residual for this data point is calculated as $60 - 60.797 = -0.797$. The negative sign confirms that the regression model slightly overestimated this individual's height by 0.797 inches.

Example 2: Calculating the Residual for Individual 2

We apply the identical systematic procedure to calculate the [residual](#) for the second individual in

our dataset.

Weight (lbs)	Height (inches)
140	60
155	62
159	67
179	70
192	71
200	72
212	75

The second individual has an observed weight (x) of **155** lbs. and an actual height (y) of **62** inches. We use the derived regression equation once again:

$$\text{height} = 32.783 + 0.2001 * (\text{weight})$$

The predicted height calculation is:

$$\text{height} = 32.783 + 0.2001 * (155)$$

$$\text{height (?) = 63.7985 inches}$$

Finally, we calculate the [residual](#) by subtracting the predicted height from the actual height:

$$\text{Residual} = 62 - 63.7985 = \mathbf{-1.7985}.$$

This larger negative residual suggests that the prediction generated by the model for the second individual was less accurate than for the first, with the model overestimating the height by nearly 1.8 inches.

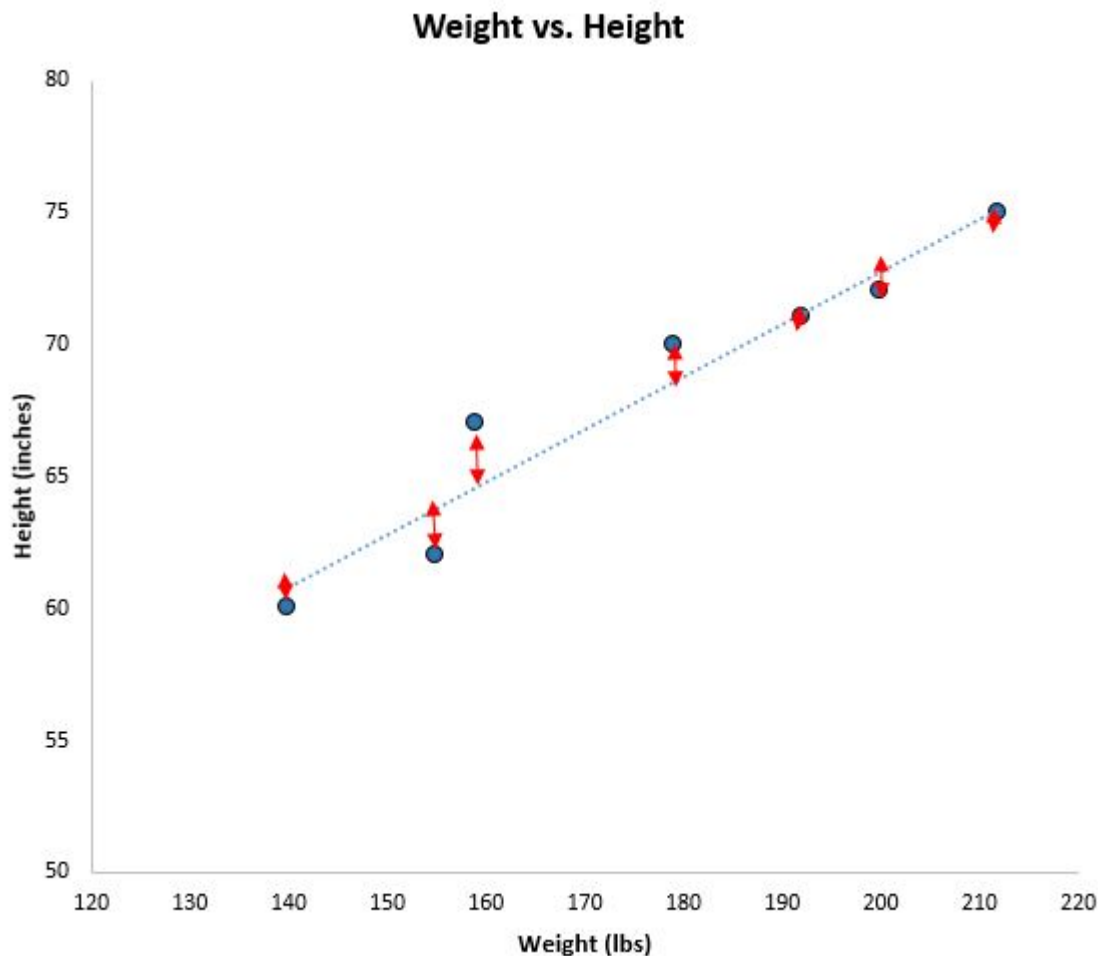
Visualizing Residuals and the Sum-to-Zero Property

By repeating the calculation process for all seven data points in our sample, we can determine the exact prediction error associated with each observation. The table below summarizes the results, clearly showing the actual values, the predicted values (?), and the corresponding [residual](#) for every individual:

Weight (lbs)	Height (inches)	Predicted Height	Residual
140	60	60.797	-0.797
155	62	63.7985	-1.7985
159	67	64.5989	2.4011
179	70	68.6009	1.3991
192	71	71.2022	-0.2022
200	72	72.803	-0.803
212	75	75.2042	-0.2042

A fundamental and essential characteristic of the [line of best fit](#), specifically when derived using the Ordinary Least Squares method, is that **if we sum up all the calculated residuals, their total will always equal zero (or be extremely close to zero due to minor rounding)**. This property is intrinsic to the OLS minimization procedure. Because the line is positioned to minimize the total squared errors, it must inherently balance the influence of all positive residuals (data points lying above the line) and all negative residuals (data points lying below the line).

A **residual** is, geometrically speaking, simply the vertical distance between the observed y-value and the value predicted by the regression line. Visualizing these vertical segments provides the clearest interpretation of the model's error for each point. The following image highlights these distances, demonstrating how some data points are far from the line (indicating large residuals and poor prediction) while others are very close (indicating small residuals and high predictive accuracy).



The Residual Plot: A Tool for Model Diagnosis

The core analytical purpose of calculating and analyzing residuals is to rigorously assess the quality and appropriateness of the [linear regression](#) line in fitting the underlying data structure. The size of the residuals serves as a direct, quantitative measure of the model's predictive accuracy. Generally, **smaller residuals** indicate that the regression line is an excellent fit, as the actual data points cluster very tightly around the predictive line. Conversely, **larger residuals** are symptomatic of a poor fit, where the model consistently deviates significantly from the actual observations.

To visualize and evaluate the entire collection of residuals simultaneously, statisticians utilize a specialized diagnostic graph known as a [residual plot](#). This plot is constructed by displaying the predicted values (?) on the horizontal axis against the corresponding calculated residual values (e) on the vertical axis. Unlike the initial scatterplot, which shows raw data, the [residual plot](#) focuses entirely on the pattern, or lack thereof, within the error component of the model.

The [residual plot](#) is indispensable for verifying several foundational assumptions of [linear regression](#), which are crucial for ensuring the statistical reliability of any inferences drawn from the

model. Specifically, analysts look for two primary ideal characteristics:

Assumption of Linearity: The residuals must be randomly and uniformly scattered around the zero horizontal line. The presence of any discernible curve or pattern (such as a U-shape or a clear trend) strongly suggests that a non-linear relationship exists between X and Y, and thus the linear model being tested is inappropriate for the data.

Constant Variance (Homoscedasticity): The vertical spread (variance) of the residuals should remain roughly constant across all predicted values displayed on the x-axis. If the variance of the residuals systematically increases or decreases as predicted values increase (often forming a distinct fan or cone shape), this condition indicates [heteroscedasticity](#), which violates a critical assumption necessary for robust linear modeling.

A truly well-fitting linear model will generate a [residual plot](#) where the points resemble pure random noise, scattered evenly and without pattern around the zero reference line. This confirms that the model has successfully captured the systematic, linear relationship and that the remaining errors are purely stochastic. For technical guidance, check out [this guide](#) to learn how to create a residual plot for a simple linear regression model in Excel.