

# How to Perform Simple Linear Regression in Excel?

Authored by  
**stats writer**

December 28, 2025

## RECOMMENDED CITATION

stats writer (2025). *How to Perform Simple Linear Regression in Excel?*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=109398>

Conducting **simple linear regression** in **Microsoft Excel** is a fundamental skill for anyone analyzing relationships between variables. The process traditionally involves several key steps: visualizing the data using a **scatter plot**, fitting a line of best fit to this visualization, and mathematically defining this relationship through a regression equation. While Excel offers specific built-in functions like **SLOPE** and **INTERCEPT** for calculating these parameters individually, the most robust and statistically complete method involves utilizing the dedicated **Data Analysis ToolPak**. This comprehensive approach provides a full statistical summary, including the equation coefficients, the ANOVA table, and the crucial **R-squared value**, which is essential for assessing the overall goodness of model fit. This thorough procedure ensures that the resulting analysis is both valid and easily interpretable, moving beyond simple calculation to complete statistical modeling.

At its core, **Simple linear regression** is a powerful statistical technique designed to model the linear relationship between two continuous variables: an independent variable (often termed the **explanatory variable**, X) and a dependent variable (the **response variable**, Y). The primary statistical objective is to find the straight line that minimizes the sum of squared errors between the predicted values and the actual observed data points. This tutorial provides a meticulous, step-by-step guide on how to execute this entire regression procedure accurately within the Excel environment, ensuring valid and high-quality statistical results suitable for academic or professional reporting.

### Prerequisite: Ensuring the Data Analysis ToolPak is Active

Before initiating any complex statistical procedure in Excel, it is essential to confirm that the **Data Analysis ToolPak** add-in is properly enabled. This tool is not activated by default in most Excel installations, yet it is indispensable for running robust analyses such as regression, ANOVA, and t-tests. Attempting to execute the regression steps without this prerequisite will result in the inability to locate the necessary **Data Analysis** option under the **Data** tab on the Excel ribbon, halting the process before it can begin.

To enable the tool, navigate through the Excel menu by clicking **File**, selecting **Options**, and then choosing **Add-ins** from the navigational pane. In the Add-ins window, locate the **Manage** dropdown menu at the bottom, select **Excel Add-ins**, and click **Go**. A subsequent dialogue box will appear, listing available add-ins. You must ensure the checkbox next to **Analysis ToolPak** is checked, and then confirm your selection by clicking **OK**. After successful activation, the **Data Analysis** button should appear visibly on the far right of the **Data** tab, granting full access to the suite of statistical tools required for the linear regression.

This activation step is paramount because standard Excel functions, while proficient for basic

mathematical calculations, do not provide the wealth of summary statistics--such as the ANOVA table, standard errors of coefficients, t-statistics, and associated P-values--that are necessary for a comprehensive, scientifically sound interpretation of a regression model. Relying on the dedicated ToolPak ensures a complete and professional output for your analysis.

## Case Study Setup: Analyzing Study Hours Versus Exam Scores

To effectively demonstrate the application of simple linear regression, we will utilize a practical and intuitive example focusing on the correlation between student effort and academic outcome. Our hypothetical scenario involves a researcher seeking to quantify the relationship between the number of hours a student dedicates to studying for an examination and the final grade achieved on that exam. The investigation aims to determine if study time is a statistically significant predictor of performance and to model the precise nature of that relationship.

In this analytical framework, **hours studied** is designated as our **explanatory variable** (X), representing the factor hypothesized to influence the outcome. Conversely, the **exam score** received acts as the **response variable** (Y), the variable whose variation we are attempting to predict and explain using the linear model. For this study, empirical data has been collected from twenty individual students, meticulously recording both their preparation time and their subsequent scores.

The initial practical requirement in Excel is the accurate transcription of this raw data into contiguous columns. Accuracy in this data entry phase is critical, as any mistakes will lead to a fundamentally biased or flawed regression model. We must designate Column A for the explanatory variable (Hours Studied) and Column B for the response variable (Exam Score), ensuring that corresponding data points--the pair of X and Y values for a single student--are perfectly aligned row by row throughout the dataset.

### Step 1: Data Entry and Preparation in Excel

#### Step 1: Enter the data.

Enter the following paired data for the number of hours studied and the exam score received for 20 students. It is paramount that the columns are properly labeled (e.g., "Hours Studied" and "Exam Score"). These labels are not just for organizational clarity; they must be included in the range selection during the regression setup, allowing the **Data Analysis ToolPak** to use meaningful names in the final output tables.

	A	B	C	D	E
1	<b>hours</b>	<b>score</b>			
2	1	76			
3	2	78			
4	2	85			
5	4	88			
6	2	72			
7	1	69			
8	5	94			
9	4	94			
10	2	88			
11	4	92			
12	4	90			
13	3	75			
14	6	96			
15	5	90			
16	3	82			
17	4	85			
18	6	99			
19	2	83			
20	1	62			
21	2	76			
22					
23					
24					

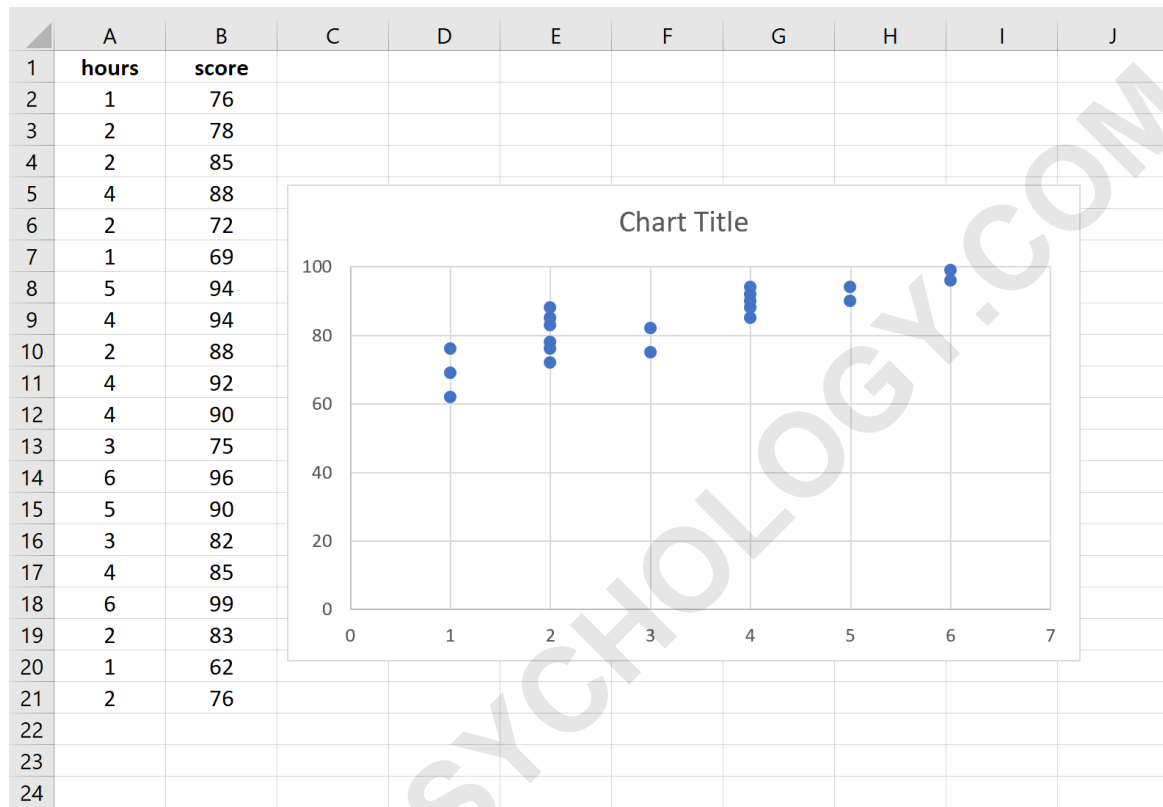
After the raw data is completely populated, a quick review of the numerical entries is prudent. While regression analysis is robust to some noise, identifying and potentially addressing extreme outliers before proceeding to the graphical representation is good statistical practice. Outliers can disproportionately influence the calculated slope and intercept of the line of best fit, potentially resulting in a misleading or inaccurate model representation of the overall population trend.

## Step 2: Visualizing the Data with a Scatter Plot

### Step 2: Visualize the data.

Before commencing the formal quantitative analysis of **simple linear regression**, the creation of a **scatter plot** of the data is a non-negotiable diagnostic step. This visual representation serves as a rapid, initial check to confirm the linearity assumption--that is, whether a straight line is indeed the most appropriate model to describe the relationship between hours studied and exam score. If the data points appear curved or randomly dispersed, a linear model would be statistically inappropriate.

To generate the plot, begin by **Highlighting** the entire data range in columns A and B, making sure to include the column headers (labels). Next, navigate to the **Insert** tab on the main Excel ribbon. Within the **Charts** group, locate and click on **Insert Scatter (X, Y)**, subsequently selecting the basic **Scatter** chart option (the first available thumbnail). Excel will automatically render the visualization of your dataset in a new chart object.



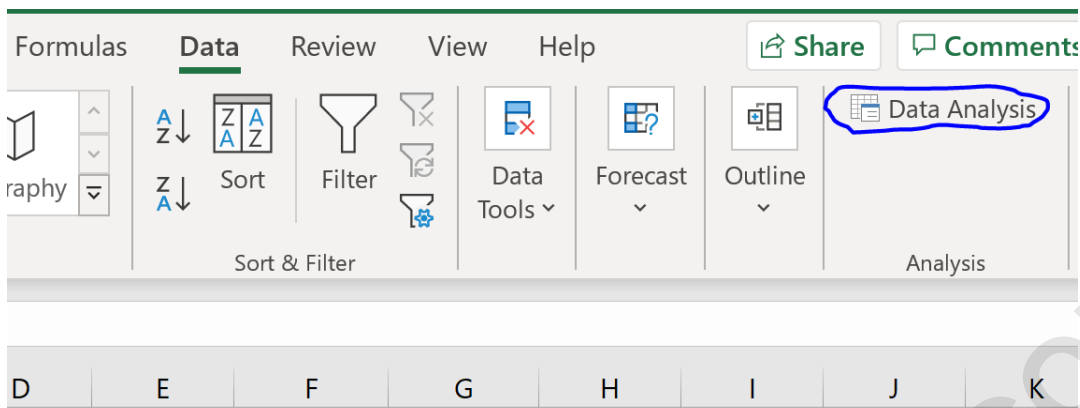
Observation of the resulting plot confirms that **hours studied** are on the horizontal (X) axis, and **exam scores** are on the vertical (Y) axis. The distribution of the data points clearly shows a strong, positive, upward-sloping pattern, indicating a robust linear relationship: greater time invested in studying is associated with higher exam scores. This strong visual evidence provides the necessary justification to proceed confidently with the formal simple linear regression quantification.

### Step 3: Utilizing the Data Analysis ToolPak for Regression

#### Step 3: Perform simple linear regression.

With the data prepared and the linearity assumption visually confirmed, the analytical computation is the next phase. Access the **Data** tab along the top ribbon in Excel and click on the **Data Analysis** option, located in the far-right section. It is a critical reminder that if this option is unavailable, the user must first **enable the Data Analysis ToolPak** add-in through the Excel

Options menu, as outlined in the prerequisite section.



Once the list of analysis tools appears, scroll through the options, select **Regression**, and then click **OK**. This action summons the specialized Regression dialogue box, requiring the precise definition of input ranges and the desired output configuration.

	A	B	C	D	E	F	G	H
1	<b>hours</b>	<b>score</b>						
2	1	76						
3	2	78						
4	2	85						
5	4	88						
6	2	72						
7	1	69						
8	5	94						
9	4	94						
10	2	88						
11	4	92						
12	4	90						
13	3	75						
14	6	96						
15	5	90						
16	3	82						
17	4	85						
18	6	99						
19	2	83						
20	1	62						
21	2	76						
22								
23								
24								

**Data Analysis** ? X

Analysis Tools

- Histogram
- Moving Average
- Random Number Generation
- Rank and Percentile
- Regression**
- Sampling
- t-Test: Paired Two Sample for Means
- t-Test: Two-Sample Assuming Equal Variances
- t-Test: Two-Sample Assuming Unequal Variances
- z-Test: Two Sample for Means

OK Cancel Help

Define the input parameters with precision:

For **Input Y Range**, select the entire array of values for the **response variable** (Exam Score),

ensuring the column label is included.

For **Input X Range**, select the entire array of values for the **explanatory variable** (Hours Studied), also including its column label.

Check the **Labels** box. This setting is mandatory when headers are included in the input ranges, instructing Excel to correctly parse the first row as identifying labels rather than data points.

For the **Output Range**, select a single empty cell (e.g., D1) on the current or a new worksheet where the extensive statistical output report should begin.

	A	B	C	D	E	F	G	H	I
1	<b>hours</b>	<b>score</b>							
2	1	76							
3	2	78							
4	2	85							
5	4	88							
6	2	72							
7	1	69							
8	5	94							
9	4	94							
10	2	88							
11	4	92							
12	4	90							
13	3	75							
14	6	96							
15	5	90							
16	3	82							
17	4	85							
18	6	99							
19	2	83							
20	1	62							
21	2	76							
22									
23									
24									
25									
26									

**Regression** dialog box settings:

- Input Y Range: \$B\$1:\$B\$21
- Input X Range: \$A\$1:\$A\$21
- Labels
- Constant is Zero
- Confidence Level: 95 %
- Output options:
  - Output Range: \$D\$2
  - New Worksheet Ply:
  - New Workbook
- Residuals:
  - Residuals
  - Residual Plots
  - Standardized Residuals
  - Line Fit Plots
- Normal Probability:
  - Normal Probability Plots

After verifying all settings, click **OK**. Excel will instantaneously generate the detailed statistical summary, ready for interpretation.

#### Step 4: Interpreting Key Regression Statistics

The resulting output is structured into three essential tables: Regression Statistics, ANOVA, and Coefficient Details. A proper conclusion to the analysis requires a systematic understanding of the key metrics presented. The following visualization represents the final summary report generated

by Excel:

D	E	F	G	H	I	J	K	L
SUMMARY OUTPUT								
<i>Regression Statistics</i>								
Multiple R	0.8528							
R Square	0.7273							
Adjusted R Square	0.7121							
Standard Error	5.2805							
Observations	20							
ANOVA								
	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>Significance F</i>			
Regression	1	1338.2906	1338.2906	47.9952	0.0000			
Residual	18	501.9094	27.8839					
Total	19	1840.2000						
	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>	<i>Lower 95.0%</i>	<i>Upper 95.0%</i>
Intercept	67.1617	2.6633	25.2178	0.0000	61.5664	72.7570	61.5664	72.7570
hours	5.2503	0.7578	6.9279	0.0000	3.6581	6.8424	3.6581	6.8424

#### Step 4: Interpret the output.

We begin by focusing on the metrics of model fit found in the Regression Statistics table:

**R Square (R<sup>2</sup>): 0.7273.** This value is the **coefficient of determination**, a measure of how well the regression line approximates the real data points. It quantifies the proportion of the total variation in the response variable (exam score) that is explained by the explanatory variable (hours studied). In this instance, 72.73% of the variation observed in the exam scores is accounted for by the number of hours a student studied. This indicates a strong predictive power, with the remaining variance attributable to unobserved confounding variables.

**Standard Error: 5.2805.** This metric represents the average deviation of the observed values from the calculated regression line, essentially serving as the standard deviation of the model's residuals. It provides a highly practical measure of the average magnitude of the error in prediction. In this case, the model predicts the exam score with an average error margin of 5.2805 points. A lower standard error signifies a model with greater predictive accuracy and a tighter fit to the underlying data.

### Evaluating Overall Model Significance: The ANOVA Table

The ANOVA (Analysis of Variance) section of the output is crucial for testing the overall statistical

reliability and significance of the entire regression model. It addresses the fundamental question: Does the explanatory variable contribute significantly to the prediction of the response variable, or is the observed relationship merely due to random sampling variability?

**F Statistic: 47.9952.** This is the calculated test statistic, which compares the variance explained by the model (Regression Mean Square) against the unexplained variance (Residual Mean Square). A substantially high F-statistic suggests that the model explains significantly more variance than would be expected by chance.

**Significance F (P-value): 0.0000.** This is the probability value (p-value) associated with the overall F statistic. It determines the probability of observing such a strong relationship if, in reality, there were no true linear association between the variables (the null hypothesis). Since this **p-value** is extremely small (approaching zero) and is far below the typical significance threshold (alpha = 0.05), we confidently reject the null hypothesis. This robust finding confirms that the overall regression model is **statistically significant**, establishing a genuine and reliable association between hours studied and exam score received.

The statistical significance provided by the ANOVA table is the foundational element that validates the use of our linear model as a true representation of the predictive relationship identified in the sample data.

## Deriving and Applying the Estimated Regression Equation

The final table detailing the coefficients provides the specific numerical parameters required to construct the estimated regression equation, which mathematically defines the line of best fit:  $Y = \text{Intercept} + (\text{Slope} * X)$ .

**Coefficients:** The values under the **Coefficients** column are the calculated intercept and slope for the relationship.

From the Excel output, the Intercept (labeled as 'Intercept') is 67.1633, and the slope coefficient (labeled as 'Hours Studied') is 5.2503. The estimated regression equation is therefore constructed as:

$$\text{exam score} = 67.16 + 5.2503 * (\text{hours studied})$$

These coefficients carry precise statistical meanings that dictate practical interpretation:

**Intercept (67.16):** Interpreted as the predicted value of the response variable when the explanatory variable is zero. In this study, 67.16 is the expected baseline exam score for a student who engages in zero hours of study.

**Slope Coefficient (5.2503):** This critical value signifies the average expected change in the exam score (Y) for every one-unit increase in hours studied (X). The interpretation is that for each **additional hour studied**, a student's exam score is expected to increase by **5.2503 points**, on average, holding all other factors constant.

This regression equation transforms into a practical predictive tool. We can use it to forecast the expected performance of any student based on their planned study time. For instance, if a student studies for exactly three hours, their expected score is calculated by substituting the value into the equation:

$$\text{expected exam score} = 67.16 + 5.2503 * (3) = 82.91$$

Thus, a student dedicating three hours to study is expected to achieve an exam score of **82.91**, offering a quantifiable prediction based on the empirical data collected.

### Advanced Diagnostic Checks for Model Validation

The successful execution and interpretation of simple linear regression in Excel provides a powerful foundation for quantitative analysis. By utilizing the **Data Analysis ToolPak**, analysts gain access to essential metrics and hypothesis tests that fully adhere to statistical standards. However, a complete analysis requires confirmation that the underlying assumptions of linear regression have been met.

To ensure the maximum reliability and trustworthiness of the model, analysts must perform diagnostic checks, primarily involving the analysis of residuals (the differences between the observed and predicted Y values). Violations of assumptions, such as non-linearity or heteroscedasticity, can render the coefficients and p-values unreliable.

Excel, in combination with its charting capabilities, can facilitate these necessary secondary analyses. Analysts often proceed with the following crucial steps to validate their models:

[How to Create a Residual Plot in Excel](#)

[How to Construct a Prediction Interval in Excel](#)

[How to Create a Q-Q Plot in Excel](#)