

How to Easily Perform Bivariate Analysis in Python

Authored by
stats writer

December 2, 2025

RECOMMENDED CITATION

stats writer (2025). *How to Easily Perform Bivariate Analysis in Python*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=103638>

Bivariate analysis is a fundamental pillar of descriptive statistics and exploratory data analysis (EDA). Unlike univariate analysis, which focuses solely on the characteristics of a single variable, bivariate analysis systematically examines the empirical relationship, association, or causality between **two distinct variables**. The core objective is to determine how changes in one variable correspond to changes in the other, offering crucial insights necessary for building predictive models or deriving meaningful statistical inferences. This process is indispensable in fields ranging from finance and social science to engineering and bioinformatics, providing the necessary foundation before moving on to more complex multivariate techniques.

The effectiveness of bivariate analysis hinges on selecting appropriate techniques based on the nature and scale of the variables involved. When analyzing two quantitative variables, common methods include generating scatter plots for visual assessment, calculating correlation coefficients for quantifying linear strength, and fitting regression models to predict outcomes. For categorical data, techniques like cross-tabulation (contingency tables) and chi-square tests are often employed. Regardless of the data type, implementing these powerful techniques is streamlined through specialized Python libraries such as Pandas for data handling, Matplotlib and Seaborn for visualization, and SciPy or Statsmodels for rigorous statistical modeling.

The Purpose and Scope of Bivariate Relationships

The primary goal when conducting bivariate analysis is to characterize the interaction between two chosen variables. This relationship can manifest in several ways: it might be **positive** (as one variable increases, the other tends to increase), **negative** (as one increases, the other tends to decrease), or there might be **no discernible linear relationship** at all. Identifying the shape and direction of this relationship is paramount because it informs subsequent analytical steps. For example, understanding a strong positive correlation between advertising spend and sales volume allows a business to prioritize investment in advertising, whereas a weak correlation might suggest diminishing returns or the presence of confounding variables that need to be explored in a multivariate context.

In the context of quantitative analysis, three techniques stand out as essential tools for assessing linear bivariate relationships, forming the backbone of many introductory statistical studies. We utilize these methods sequentially: first, visual inspection via a scatter plot provides an intuitive understanding of the data distribution; second, calculating a correlation coefficient offers a precise, standardized numerical measure of the association strength; and finally, simple linear regression provides a predictive equation, allowing us to quantify the exact impact of the explanatory variable on the response variable. Mastery of these three approaches ensures a comprehensive analysis of the interaction between any pair of continuous data points within a dataset.

Setting Up the Python Environment and Data

Before diving into the analytical methods, it is standard procedure to prepare and structure the data using the [Pandas](#) library, the cornerstone of data manipulation in [Python](#). For this illustrative example, we will analyze a synthetic dataset designed to explore the relationship between the number of hours a student spends studying and the corresponding exam score received. This dataset contains observations for 20 distinct students, allowing us to treat 'hours studied' as the explanatory variable (X) and 'exam score' as the response variable (Y). Initial data preparation involves importing the necessary library and constructing the Pandas DataFrame, ensuring the data is correctly indexed and typed before any statistical operations commence.

The following code snippet demonstrates the creation of this DataFrame. Notice how the structure clearly defines two columns: [Pandas DataFrames](#) are optimized for columnar data, making operations like correlation calculation or plotting incredibly efficient. Reviewing the initial rows using the `.head()` method is a vital step in data quality assurance, confirming that the data was loaded correctly and matches expectations. This foundational step ensures that subsequent analytical processes are robust and reliable.

```
import pandas as pd
```

```
#create DataFrame
```

```
df = pd.DataFrame({'hours': ,  
'score': })
```

```
#view first five rows of DataFrame
```

```
df.head()
```

```
hours score
```

```
0 1 75
```

```
1 1 66
```

```
2 1 68
```

```
3 2 74
```

```
4 2 78
```

Method 1: Visualizing Relationships with Scatter Plots

The most intuitive and crucial first step in any [bivariate analysis](#) is graphical visualization, typically achieved using a [scatter plot](#). A scatter plot maps pairs of data points (X, Y) onto a Cartesian plane, where the independent variable (Hours Studied) is placed on the horizontal X-axis and the dependent variable (Exam Score) is placed on the vertical Y-axis. This visualization is essential because it allows the analyst to immediately assess three key characteristics of the relationship:

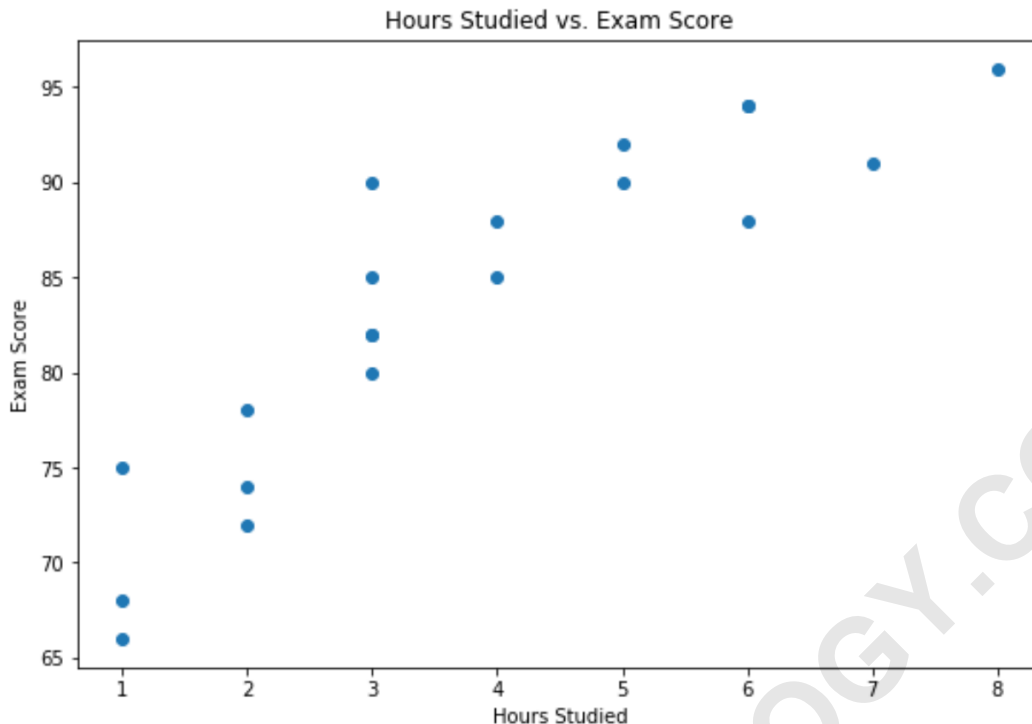
the **form** (linear, curvilinear, or none), the **direction** (positive or negative), and the **strength** (how closely the points cluster around a potential line).

We utilize the powerful Matplotlib library, specifically its `pyplot` module, to generate this visualization in Python. The syntax below demonstrates how to call the `.scatter()` function, passing the relevant columns from our Pandas DataFrame. Furthermore, proper labeling of the axes using `.xlabel()` and `.ylabel()` and defining a clear title using `.title()` are critical steps in creating a professional and easily interpretable figure, adhering to best practices in data communication.

import matplotlib.pyplot as plt

```
#create scatterplot of hours vs. score
plt.scatter(df.hours, df.score)
plt.title('Hours Studied vs. Exam Score')
plt.xlabel('Hours Studied')
plt.ylabel('Exam Score')
```

Upon reviewing the generated plot, we can visually confirm the nature of the relationship between the variables. The clustering of data points, as illustrated in the figure below, shows a distinct upward trend from left to right. This pattern unequivocally indicates a **strong positive relationship**: students who dedicated more hours to studying generally achieved higher exam scores. This visual evidence serves as a crucial confirmation before proceeding to calculate numerical measures of association, ensuring that our statistical models align with the observable data pattern.



Method 2: Quantifying Linear Association using Correlation Coefficients

While a scatter plot provides qualitative insight into the relationship, a correlation coefficient offers a standardized, quantitative measure of the strength and direction of the linear association between two continuous variables. The most commonly used metric for this purpose is the Pearson Correlation Coefficient (r), which always falls between -1 and +1. A value of +1 signifies a perfect positive linear relationship, -1 indicates a perfect negative linear relationship, and 0 indicates no linear correlation whatsoever. It is vital to remember that correlation only measures linear association and does not imply causation between the two variables.

Calculating the correlation coefficient in Python is straightforward thanks to the optimized functions available within the Pandas library. By simply calling the `.corr()` method on the DataFrame, Pandas generates a **correlation matrix**. This matrix displays the correlation coefficient for every possible pairing of variables within the DataFrame. Since our DataFrame only contains 'hours' and 'score', the resulting 2x2 matrix provides the correlation of each variable with itself (which is always 1) and the correlation between the two variables of interest.

```
#create correlation matrix
```

```
df.corr()
```

```
hours score
```

```
hours 1.000000 0.891306
```

score 0.891306 1.000000

Analyzing the output matrix reveals that the correlation coefficient between 'hours' and 'score' is approximately **0.891**. This numerical result strongly confirms the initial visual assessment from the scatter plot. A coefficient this close to +1 is conventionally categorized as a **very strong positive correlation**, suggesting that increased study time is highly reliably associated with increased exam performance within this sample population. This high correlation coefficient justifies proceeding to the next step: using simple linear regression to build a predictive model.

Method 3: Modeling Relationships with Simple Linear Regression

The final and most advanced step in this bivariate analysis is applying Simple Linear Regression (SLR). SLR moves beyond mere association (correlation) and attempts to model the relationship, allowing us to predict the value of the dependent variable (Y) based on the value of the independent variable (X). The model seeks to find the "line of best fit"--the line that minimizes the sum of squared residuals--following the mathematical form $Y = \beta_0 + \beta_1 X + \epsilon$, where β_0 is the intercept, β_1 is the slope coefficient, and ϵ represents the error term.

To perform SLR in Python, we typically leverage the `statsmodels` library, which provides comprehensive statistical outputs comparable to specialized statistical software. We define our response (Y, score) and explanatory (X, hours) variables. A critical step when using `statsmodels` for Ordinary Least Squares (OLS) estimation is adding a constant term to the explanatory variables using `sm.add_constant(x)`. This constant represents the intercept (β_0) in the regression equation and must be explicitly included for the model to calculate it correctly.

The core of the analysis involves fitting the model using `sm.OLS(y, x).fit()`. This command executes the least squares estimation process. The resulting object, `model`, contains all the necessary statistical information, which can then be printed using `model.summary()`. This summary is rich with details, including coefficient estimates, standard errors, test statistics, p-values, and overall model performance metrics like R^2 .

```
import statsmodels.api as sm
```

```
#define response variable
```

```
y = df
```

```
#define explanatory variable
```

```
x = df]
```

```
#add constant to predictor variables
```

```
x = sm.add_constant(x)

#fit linear regression model
model = sm.OLS(y, x).fit()

#view model summary
print(model.summary())
```

OLS Regression Results

```
=====
===
Dep. Variable: score R-squared: 0.794
Model: OLS Adj. R-squared: 0.783
Method: Least Squares F-statistic: 69.56
Date: Mon, 22 Nov 2021 Prob (F-statistic): 1.35e-07
Time: 16:15:52 Log-Likelihood: -55.886
No. Observations: 20 AIC: 115.8
Df Residuals: 18 BIC: 117.8
Df Model: 1
Covariance Type: nonrobust
=====
===
coef std err t P>|t|
-----
const 69.0734 1.965 35.149 0.000 64.945 73.202
hours 3.8471 0.461 8.340 0.000 2.878 4.816
=====
===
Omnibus: 0.171 Durbin-Watson: 1.404
Prob(Omnibus): 0.918 Jarque-Bera (JB): 0.177
Skew: 0.165 Prob(JB): 0.915
Kurtosis: 2.679 Cond. No. 9.37
=====
===
```

Interpreting the Regression Output

The summary table generated by `statsmodels` provides all the necessary parameters for interpreting the fitted model. The most critical information is contained within the middle coefficient table. From this output, we extract the estimated coefficients, which define our predictive model.

The coefficient for the `const` (intercept) is \$69.0734\$, and the coefficient for `hours` (the slope, `$beta_1$`) is \$3.8471\$. These values allow us to formulate the definitive regression equation for predicting exam scores based on study hours:

Predicted Exam Score = 69.0734 + 3.8471 * (Hours Studied)

The interpretation of these coefficients is statistical: The intercept (\$69.0734\$) represents the predicted exam score when the hours studied is zero. The slope coefficient (\$3.8471\$) signifies that, on average, for every **one-unit increase** in hours studied, the predicted exam score increases by 3.8471 points. Furthermore, we examine the R^2 value, which is 0.794. This means that 79.4% of the variability observed in the exam scores can be explained by the variation in the hours studied, confirming that hours studied is a highly significant predictor variable in this bivariate relationship.

Finally, the simple linear regression model is used for prediction. We can input any value for hours studied (within the observed range) into the derived equation to estimate the corresponding exam score. For instance, if a new student reported studying for 3 hours, their predicted score would be calculated precisely, demonstrating the practical application of our statistical model in generating actionable predictions.

The calculation for predicting a student score after 3 hours of study is as follows:

Exam Score = 69.0734 + 3.8471 * (hours studied)

Exam Score = 69.0734 + 3.8471 * (3)

Exam Score = 81.6147

Summary and Advanced Considerations

Mastering bivariate analysis is critical for any data science professional, forming the bridge between raw data observation and complex statistical modeling. By systematically employing visualization techniques like the scatter plot, quantifying the linear strength using the Pearson Correlation Coefficient, and deriving a predictive model through Simple Linear Regression, we gain a complete understanding of how two variables interact. In our study of student performance, we established a strong, positive, and statistically significant relationship between study hours and exam scores.

While this article focused on continuous variables and linear relationships using Python libraries like Pandas, Matplotlib, and Statsmodels, it is important to acknowledge the broader scope of bivariate techniques. Analysts must also be prepared to handle non-linear associations, which might require transformations or alternative models, and categorical data analysis, which involves methods such as Analysis of Variance (ANOVA) or chi-squared tests. Continuous scrutiny of

assumptions--such as linearity, homoscedasticity, and normality of residuals in regression--is always necessary to ensure the validity and reliability of the statistical inferences drawn from the analysis.

ARABPSYCHOLOGY.COM