

How to Understand Causation vs. Correlation: 3 Simple Examples

Authored by
stats writer

December 2, 2025

RECOMMENDED CITATION

stats writer (2025). *How to Understand Causation vs. Correlation: 3 Simple Examples*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=103839>

The relationship between causation and correlation is one of the most fundamental, yet frequently misunderstood, concepts in statistics and scientific inquiry. While often confused, these two principles describe distinct phenomena. Correlation simply measures the degree to which two variables move together, indicating a statistical association. Causation, however, describes a definitive mechanism where a change in one variable (the cause) directly produces a change in another variable (the effect). It is a well-established statistical truth that correlation does not imply causation, preventing us from making definitive claims based solely on shared trends. However, a deeper and equally crucial question arises when we reverse the premise: If we definitively establish that A causes B, must A and B necessarily be correlated?

The short, definitive answer is: Causation does not automatically imply correlation, especially when utilizing standard linear measures of association. This unexpected outcome stems from the limitations of the most commonly used statistical tool for measuring association, the linear correlation coefficient. When the causal relationship is highly non-linear, symmetrical, or complex, the standard linear correlation metric can misleadingly register a value of zero, thereby suggesting no association exists, even though a clear, deterministic cause-and-effect link is present. Understanding this statistical nuance is essential for researchers and analysts who rely on statistical modeling to draw accurate conclusions about complex systems.

The Critical Distinction: Causation vs. Correlation

Before exploring the main thesis, it is beneficial to firmly distinguish between these two concepts, as the frequent mixing of the terms leads to significant analytical errors. Correlation is a quantitative measure that details the strength and direction of a linear relationship between two variables. If Variable X increases and Variable Y tends to increase, they are positively correlated. If X increases and Y tends to decrease, they are negatively correlated. If they vary independently, they show zero correlation. This calculation, typically represented by the Pearson correlation coefficient (r), ranges from -1 to +1, describing only the linear alignment of the data points.

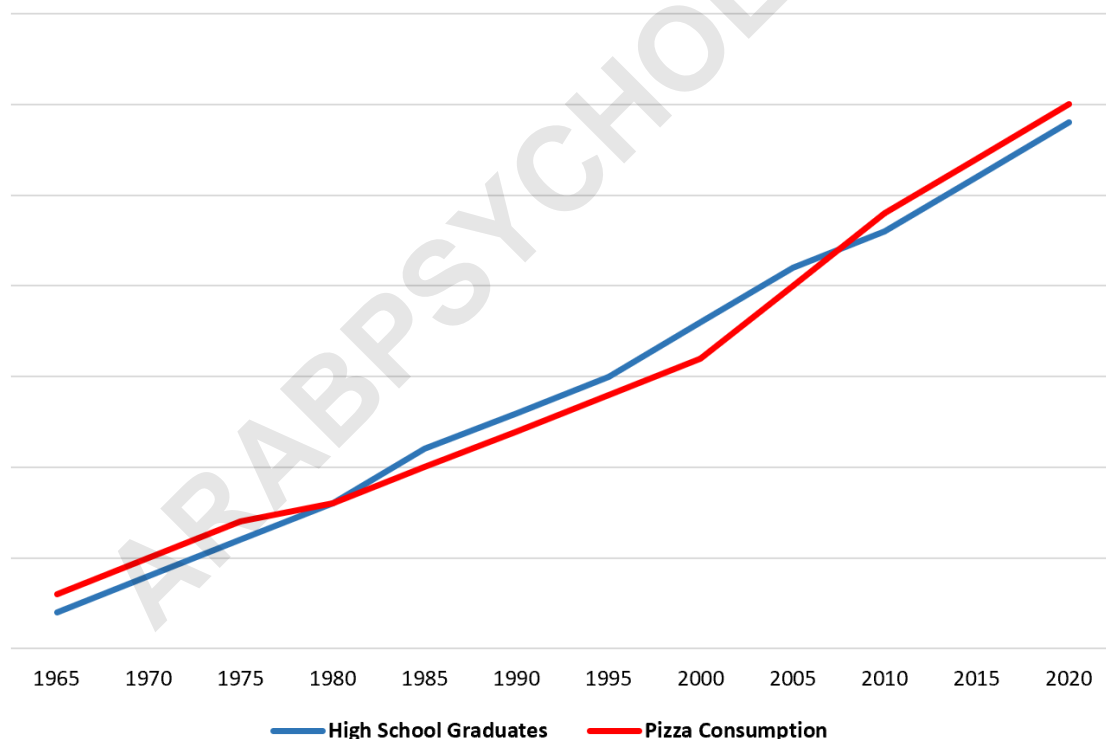
In contrast, causation requires evidence of a mechanism, temporal precedence, and the exclusion of alternative explanations. Establishing causation often requires highly controlled experimental settings, such as randomized controlled trials, or advanced statistical techniques designed to mitigate the influence of confounding variables. It is important to note that while causation is a much stronger claim than correlation, standard statistical analysis frequently focuses on correlation because it is mathematically simpler to compute and interpret, especially in preliminary data exploration. The statistical community widely accepts the principle that a strong correlation does not provide sufficient evidence to conclude causation.

The Fallacy of Correlation Implying Causation (Review)

The classic statistical warning--correlation is not causation--highlights the risk of assuming a cause-and-effect link simply because two variables show a shared trend. This often occurs when a third, unobserved variable (a confounding variable) drives both variables simultaneously. The common example used below perfectly illustrates how two seemingly related trends can be driven by a broader, underlying factor without any direct connection between them. Recognizing this allows us to set the framework for examining the reverse relationship: whether an established causal link requires a correlated output.

Consider the highly cited, albeit spurious, example comparing the number of high school graduates and the total pizza consumption in the U.S. over several decades. Both data sets exhibit a clear, positive association, meaning that as the number of graduates increases, so does the amount of pizza consumed. If one were to calculate the linear correlation, it would likely be very high. This strong statistical link, however, is merely coincidental, not causal.

High School Graduates vs. Pizza Consumption



This image demonstrates a powerful statistical association between the two variables. Despite the visually compelling trend, this does not mean that increasing the number of high school graduates is the primary factor **causing** an increase in pizza consumption. The more rational and accurate explanation is that the overall U.S. population has steadily increased over the same period. Since the population is rising, the number of people achieving a high school degree and the total volume

of food, including pizza, being consumed are both increasing as a function of the rising population. The rising U.S. population acts as the confounding variable that drives both trends independently, resulting in a spurious correlation.

Does Causation Necessarily Imply Correlation? The Statistical Reality

Having established the dangers of confusing correlation for causation, we pivot to the central question of this analysis: If one variable, X , is definitively known to cause changes in another variable, Y , does this causal relationship guarantee a measurable, non-zero linear correlation? This question is vital because many researchers rely solely on the correlation coefficient (r) as an initial screen to determine if any relationship, linear or otherwise, exists between variables. If causation is present but correlation is zero, this initial screening step fails spectacularly.

The key to understanding the answer lies in the mathematical definition of linear correlation, such as the Pearson correlation coefficient. This statistic is designed explicitly to measure the extent to which the variables adhere to a straight-line relationship. It is inherently insensitive to certain types of non-linear patterns, especially those that are symmetrical around a central point. When a causal relationship follows a curve that exhibits symmetry--such as a parabola or a cosine wave--the positive and negative deviations from the mean tend to cancel each other out when the correlation is calculated across a sufficiently wide domain, resulting in a correlation close or equal to zero. Thus, the existence of a definitive causal link is not sufficient to guarantee a measurable linear correlation.

Example 1: The Quadratic Relationship ($Y = X^2$)

Consider a variable X that deterministically causes Variable Y according to a quadratic function, where Y is simply the square of X . This is a clear, non-probabilistic causal relationship: for every input of X , the output Y is mathematically fixed. If we sample X across a symmetrical range--for instance, from -10 to $+10$ --we observe a highly predictable, parabolic pattern. The definition of a quadratic relationship ensures that X is causing Y , yet when we subject this data set to a standard linear correlation coefficient calculation, the result is zero.

The reason for this zero correlation lies in the symmetry. When X is negative (e.g., -5), Y is positive (25). As X moves toward zero, Y decreases. When X becomes positive (e.g., $+5$), Y is also positive (25), and as X increases further, Y increases again. The linear correlation algorithm averages the products of the deviations from the mean. Since the negative values of X are associated with increasing Y values, and the positive values of X are also associated with increasing Y values, the overall linear trend is neutralized. There is no consistent slope (positive or negative) that defines the relationship across the entire range, hence the zero correlation.

Suppose some variable, X , causes variable Y to take on a value equal to X^2 .

For example:

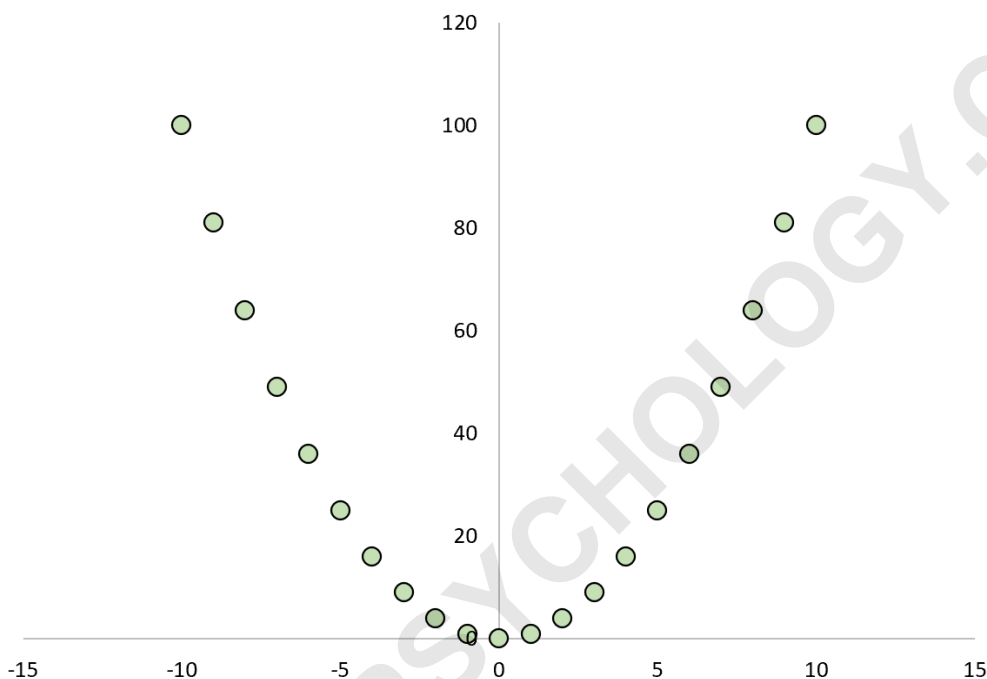
If $X = -10$ then $Y = -10^2 = 100$

If $X = 0$ then $Y = 0^2 = 0$

If $X = 10$ then $Y = 10^2 = 100$

And so on.

If we plotted the relationship between X and Y, it would look like this:



If we calculated the Pearson correlation coefficient between the two variables, we would find that the correlation is **zero**. This demonstrates a clear case where a perfect causal relationship yields zero linear correlation.

Example 2: The Quartic Relationship ($Y = X^4$)

To further illustrate the insensitivity of linear correlation to symmetrical causal patterns, we can examine a quartic relationship. If Variable X causes Variable Y such that Y is equal to X raised to the power of four, the causal link is undeniable. However, the resulting plotted relationship is even flatter and wider at the base than the quadratic relationship, further emphasizing the non-linear structure. The high degree of symmetry around $X=0$ again ensures that any calculation of linear association will fail to capture the underlying dependency.

In this quartic function, as X increases or decreases from zero, Y always increases very rapidly.

The fact that negative X values produce large positive Y values, just as positive X values do, ensures that there is no consistent positive or negative slope defining the relationship across the domain. The correlation metric interprets this lack of consistent linear direction as an absence of association. This scenario is a powerful demonstration that even highly deterministic relationships based on simple mathematical rules can remain invisible if only tested using linear statistical methods. This highlights the necessity of visual inspection and the application of appropriate Non-linear relationships modeling techniques.

Suppose some variable, X, causes variable Y to take on a value equal to X^4 ;

For example:

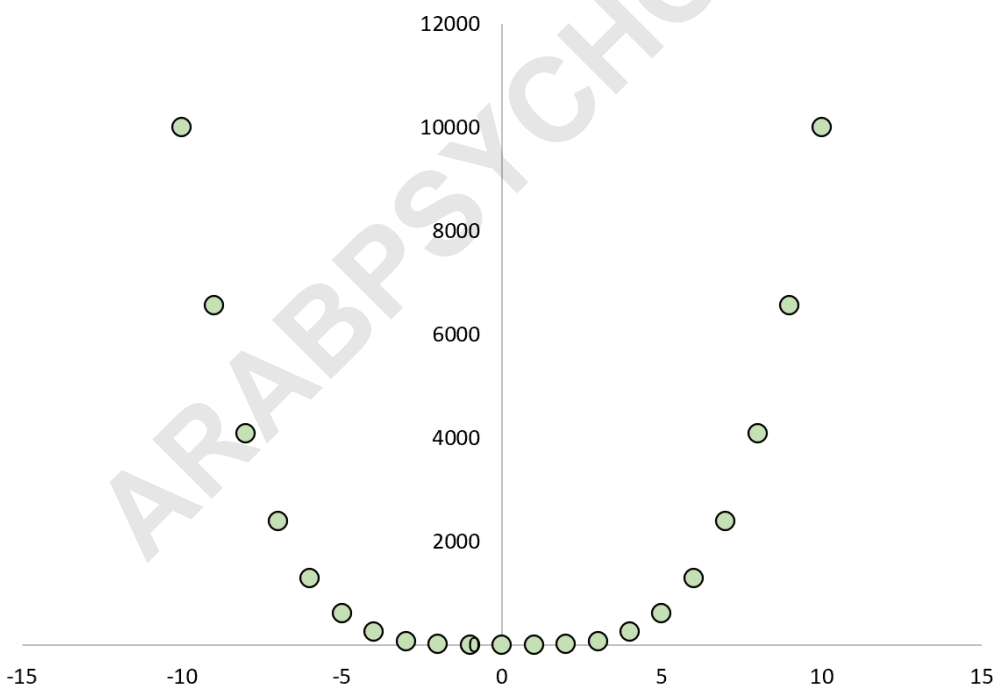
If $X = -10$ then $Y = -10^4 = 10,000$

If $X = 0$ then $Y = 0^4 = 0$

If $X = 10$ then $Y = 10^4 = 10,000$

And so on.

If we plotted the relationship between X and Y, it would look like this:



If we calculated the Pearson correlation coefficient between the two variables, we would find that the correlation is **zero**. We know that X causes Y, but the linear correlation between the two variables is zero.

Example 3: The Cosine Relationship ($Y = \cos(X)$)

The principles of symmetry and non-linearity leading to zero correlation extend beyond polynomial functions to periodic functions like the cosine wave. A periodic function ensures that the variables are causally linked, as the value of Y is entirely dependent on the value of X . If we analyze the cosine function over several periods, the relationship between X and Y becomes cyclical, not linear. As X increases, Y oscillates between -1 and $+1$.

When calculating the linear correlation over an interval that covers multiple full cycles of the cosine wave, the positive slopes are perfectly counterbalanced by the negative slopes. For instance, in some parts of the curve, as X increases, Y decreases (negative slope), while in other parts, as X increases, Y increases (positive slope). Since the Pearson correlation coefficient summarizes the overall linear tendency across the dataset, these opposing slopes cancel out, resulting in a zero correlation despite the mathematically guaranteed causation. This third example provides conclusive evidence that deterministic causation does not necessitate a linear association.

Suppose some variable, X , causes variable Y to take on a value equal to $\cos(X)$.

For example:

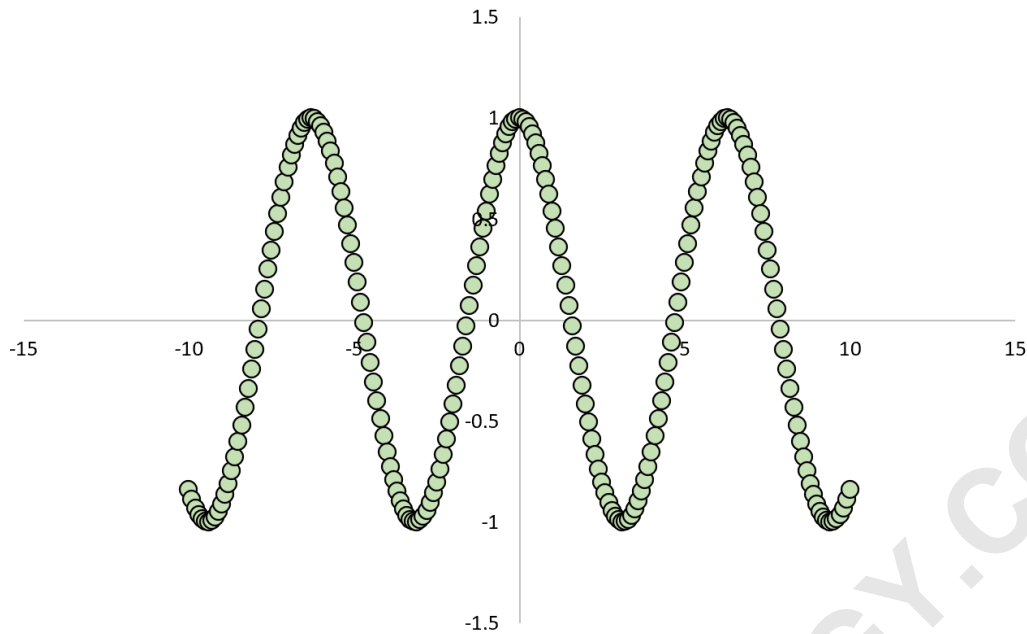
If $X = -10$ then $Y = \cos(-10) = -0.83907$

If $X = 0$ then $Y = \cos(0) = 1$

If $X = 10$ then $Y = \cos(10) = -0.83907$

And so on.

If we plotted the relationship between X and Y , it would look like this:



If we calculated the Pearson correlation coefficient between the two variables, we would find that the correlation is **zero**. We know that X causes Y, but the linear correlation between the two variables is zero.

The Importance of Testing for Non-linear relationships

The statistical truth that causation does not imply correlation has profound practical implications for researchers. If a scientist relies solely on a preliminary correlation check (e.g., calculating Pearson's r) to decide whether to pursue a causal hypothesis, they risk discarding valid and potentially important causal relationships simply because those relationships are non-linear or symmetrical. In fields ranging from physics to economics, where many underlying mechanisms follow complex, Non-linear relationships, assuming a linear form is a dangerous oversimplification.

To avoid this error, researchers must employ techniques beyond simple linear analysis. Visualizing the data through scatter plots is crucial, as the symmetrical shapes shown in the examples above are immediately apparent graphically, even if the linear correlation is zero. Furthermore, methods that test for non-linear dependence, such as polynomial regression, mutual information, or distance correlation, must be employed. These advanced metrics can capture the full extent of the statistical dependency, confirming that even when linear correlation is absent, the variables are far from statistically independent. Therefore, while correlation is a necessary condition for causation in simple linear cases, it is not a universally necessary condition for causation in the broader landscape of complex phenomena.

The following tutorials provide additional information about correlation and causation: