

How to Calculate Polychoric Correlation in R with the 'polychor' Package

Authored by
stats writer

December 3, 2025

RECOMMENDED CITATION

stats writer (2025). *How to Calculate Polychoric Correlation in R with the 'polychor' Package*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=104380>

The calculation of correlation coefficients is a fundamental task in statistical analysis, yet the choice of method is highly dependent on the nature of the data. When dealing with ordinal variables--data where categories possess a natural, meaningful order but the intervals between them are undefined--standard measures like Pearson correlation can yield misleading results. This is where the Polychoric correlation (PCC) becomes indispensable.

In the statistical programming environment R, the polychoric correlation can be robustly calculated using specialized packages, most notably the polycor package. This package implements a sophisticated statistical procedure, often relying on Maximum Likelihood Estimation (MLE), to accurately determine the association. The essence of this technique lies in modeling the relationship between the observed ordinal data based on an underlying continuous structure, typically involving unobserved latent variables.

By accounting for this hypothesized continuous distribution, the polychoric correlation provides a far more accurate and theoretically sound estimate of the true relationship between two ordinal variables than methods that treat the categories merely as numerical ranks. This depth of estimation is crucial for high-stakes research in fields ranging from psychology and sociology to market research and econometrics.

The Polychoric Correlation: Bridging the Gap Between Ordinal and Continuous Data

The Polychoric correlation coefficient is specifically designed to estimate the correlation between two theoretical continuous variables that are assumed to underlie two observed discrete, ordered variables. It is the correlation that would be observed if the measured ordinal variables were instead measured on a perfect continuous scale. This measure is essential because ordinal scaling often compresses underlying continuous traits (e.g., measuring intelligence or satisfaction on a 5-point scale rather than an infinite range).

The application of PCC is standard practice when analyzing datasets derived from surveys or assessments where responses are collected using Likert scales or other graded categories. Traditional methods fail to accurately capture the relationship because they either assume equal intervals between ranks (treating them as interval data) or focus solely on the rank order (ignoring potential parametric properties). PCC overcomes these limitations by modeling the thresholds that divide the underlying continuous distribution into the observed categories.

Like the Pearson correlation coefficient, the value for polychoric correlation ranges strictly from **-1** to **1**. This range allows for easy interpretation regarding the strength and direction of the linear association between the two underlying continuous constructs.

-1 indicates a perfect negative correlation, meaning that as one latent variable increases, the other decreases consistently.

0 indicates no correlation or linear association between the two underlying variables.

1 indicates a perfect positive correlation, meaning the two latent variables increase or decrease together perfectly.

Understanding Ordinal Variables: The Foundation of PCC Analysis

Before diving into calculation, it is vital to firmly grasp what constitutes an ordinal variable. These are variables whose possible values are categorical, but unlike nominal data, they possess a clear, intrinsic, and natural order. However, the critical distinction is that the distance between the categories (the intervals) cannot be assumed to be equal or meaningful. For instance, the difference in satisfaction between "Very Unsatisfied" and "Unsatisfied" might not be the same magnitude as the difference between "Satisfied" and "Very Satisfied."

Recognizing ordinal variables prevents the misuse of inappropriate statistical techniques. If these variables were mistakenly treated as interval or ratio data, the resulting statistical inferences--including correlation coefficients--would be biased and potentially invalid. The choice of PCC is a direct acknowledgment of this measurement limitation, seeking to estimate the continuous truth hidden beneath the ordered categories.

Several common variables used in social sciences and market research are measured on an ordinal scale. These examples illustrate the common contexts where polychoric correlation is required:

Satisfaction Assessment: Categories such as Very unsatisfied, unsatisfied, neutral, satisfied, very satisfied.

Socioeconomic Status: Groupings like Low income, medium income, high income.

Professional Hierarchy: Ranks such as Entry Analyst, Analyst I, Analyst II, Lead Analyst.

Clinical Metrics: Measurements of pain severity like Small amount, medium amount, high amount.

The Underlying Theory: Latent Variables and Maximum Likelihood Estimation

The power of the Polychoric correlation stems from its foundation in structural equation modeling principles, specifically the concept of latent variables. A latent variable is a hypothetical construct that is not directly observed but is inferred from other observable variables. In the context of PCC, we hypothesize that each ordinal variable we observe is merely a categorized manifestation of a deeper, continuous, normally distributed latent variable.

The observed ordinal categories act as thresholds slicing across this continuous latent distribution. The PCC calculation, therefore, involves estimating the correlation between these two hypothetical

continuous variables (the latent variables) that, when categorized according to certain threshold parameters, would best reproduce the observed joint frequency distribution (cross-tabulation) of the two ordinal variables.

This estimation process typically employs Maximum Likelihood Estimation (MLE). MLE is an advanced statistical method that seeks to find the parameter values (in this case, the latent correlation and the categorization thresholds) that maximize the probability (likelihood) of observing the actual data collected. Because MLE provides asymptotically efficient and unbiased estimates under specific distributional assumptions (usually bivariate normality of the latent variables), it is the standard and most robust approach for calculating PCC.

Implementing Polychoric Correlation in R: Utilizing the **polychor** Package

The statistical programming language **R** offers excellent functionality for calculating specialized correlation coefficients. To compute the polychoric correlation, researchers typically rely on the external package **polychor**, which provides the main function `polychor()`. This function is designed to efficiently handle the complex iterative calculations required by the MLE procedure.

Before executing the analysis, ensure the **polychor** package is installed in your R environment. If it is not already installed, you would use `install.packages("polychor")`. Once installed, the package must be loaded into the current session using the `library(polychor)` command. The primary syntax for the calculation is straightforward: `polychor(x, y)`, where `x` and `y` are the vectors containing the values of the two ordinal variables.

It is important to note that while the input variables `x` and `y` are numerical representations of the ordered categories (e.g., 1, 2, 3), the `polychor()` function treats them appropriately as ordinal data, performing the necessary underlying calculations based on the assumed continuous structure. This abstraction simplifies the user experience while maintaining statistical rigor, enabling immediate and accurate estimation of the underlying relationship.

Example 1: Calculating Polychoric Correlation for Movie Ratings Consistency

Imagine a scenario where we are assessing the consistency between two independent movie rating agencies. We hypothesize that if the agencies are measuring the same latent construct (film quality), their ratings, even though they are discrete ordinal scores, should exhibit a high degree of correlation. To test this, we ask each agency to rate 20 different movies on a standardized scale of 1 to 3.

The rating scale used is explicitly ordinal:

1 indicates "bad" quality (lowest rank).

2 indicates "mediocre" quality.

3 indicates "good" quality (highest rank).

We need to utilize R and the `polychor()` function to calculate the polychoric correlation between the ratings provided by Agency 1 and Agency 2. This will reveal the estimated correlation between their underlying, continuous latent evaluations of film quality.

The following R code demonstrates the setup of the data vectors and the execution of the correlation calculation. Notice how we first load the **polychor** library and then define the vectors representing the 20 ratings from each agency:

library(polychor)

```
#define movie ratings for each agency
agency1 <- c(1, 1, 2, 2, 3, 2, 2, 3, 2, 3, 3, 2, 1, 2, 2, 1, 1, 1, 2, 2)
agency2 <- c(1, 1, 2, 1, 3, 3, 3, 2, 2, 3, 3, 3, 2, 2, 2, 1, 2, 1, 3, 3)

#calculate polychoric correlation between ratings
polychor(agency1, agency2)
```

0.7828328

Upon execution, the resulting Polychoric correlation estimate is approximately **0.78**. This robust positive value indicates a strong, substantive association between the ratings assigned by the two agencies, suggesting high reliability and agreement in their underlying assessment criteria.

Interpreting the Polychoric Correlation Coefficient

Interpreting the polychoric coefficient follows the same conventions as interpreting Pearson's r , focusing on the sign (direction) and magnitude (strength). A value of **0.78**, as calculated in Example 1, is typically considered a high correlation in social science research. This means that as Agency 1 rates a movie higher on its latent quality scale, Agency 2 is also highly likely to rate that same movie higher on its latent scale.

It is paramount to remember what this coefficient represents: the correlation between the **latent continuous variables**, not the observed ordinal scores themselves. This distinction is critical because if we had mistakenly used the standard Pearson correlation on these ordinal scores, the resulting coefficient would likely be attenuated (closer to zero) or otherwise inaccurate, providing a misleading assessment of the agencies' true agreement.

A correlation near 1 implies near-perfect agreement in the underlying continuous constructs being measured. Conversely, a correlation near -1 indicates an inverse relationship, where high scores

on one construct correspond systematically to low scores on the other. A value close to 0 suggests the underlying continuous variables are independent of one another, indicating that the agencies' rating behaviors are unrelated.

Example 2: Analyzing Polychoric Correlation for Restaurant Customer Satisfaction

In our second practical example, we explore customer satisfaction ratings between two competing restaurants. We hypothesize that customers who enjoy one restaurant might also enjoy the other, or perhaps they have opposite preferences. To analyze this, we randomly surveyed 20 customers who had recently dined at both Restaurant 1 and Restaurant 2. We asked them to rate their overall satisfaction on a standard 5-point Likert scale.

The 5-point ordinal scale provided the following categories:

- 1 indicates "very unsatisfied" (lowest satisfaction).
- 2 indicates "unsatisfied".
- 3 indicates "neutral".
- 4 indicates "satisfied".
- 5 indicates "very satisfied" (highest satisfaction).

The goal is to calculate the polychoric correlation to estimate the association between the underlying continuous satisfaction level a customer holds for Restaurant 1 versus Restaurant 2. This will tell us if a customer's inherent predisposition towards liking or disliking dining experiences applies similarly across both establishments.

We again use the `polychor()` function from the **polychor** package in R. The data vectors below represent the raw ratings collected from the 20 customers for both restaurants:

library(polychor)

```
#define ratings for each restaurant
```

```
restaurant1 <- c(1, 1, 2, 2, 2, 3, 3, 3, 2, 2, 3, 4, 4, 5, 5, 4, 3, 4, 5, 5)
```

```
restaurant2 <- c(4, 3, 3, 4, 3, 3, 4, 5, 4, 4, 4, 5, 5, 4, 2, 1, 1, 2, 1, 4)
```

```
#calculate polychoric correlation between ratings
```

```
polychor(restaurant1, restaurant2)
```

```
-0.1322774
```

The resulting Polychoric correlation coefficient is calculated as approximately **-0.13**. This value is

exceptionally close to zero, which necessitates a conclusion of minimal or negligible association between the underlying customer satisfaction levels for the two restaurants. In practical terms, a customer's latent satisfaction with Restaurant 1 does not predict their latent satisfaction with Restaurant 2.

Why Polychoric Correlation Outperforms Simpler Methods for Ordinal Data

When faced with ordinal variables, analysts often debate whether to use non-parametric measures like Spearman's Rho or Kendall's Tau, or simply to treat the variables as interval data and use Pearson's r . The Polychoric correlation method provides a compelling advantage over these alternatives, especially when the goal is to estimate the relationship between the theoretical continuous traits that define the measures.

Pearson's r is fundamentally inappropriate for ordinal data because it assumes interval properties--that the distance between 1 and 2 is the same as 4 and 5--which is false for Likert scales. Using Pearson's r on ordinal data typically leads to an underestimation of the true underlying correlation. Conversely, non-parametric measures (Spearman's Rho and Kendall's Tau) focus only on the rank order. While valid, they do not utilize the assumption of an underlying continuous bivariate normal distribution, which is often a reasonable and powerful assumption in psychological and social research.

By leveraging the assumption of underlying continuous latent variables and utilizing sophisticated Maximum Likelihood Estimation, PCC offers an estimate of the correlation that is statistically more efficient and theoretically accurate under the conditions where the ordinal data truly represent categorized continuous scores. Thus, for researchers focused on modeling latent traits, PCC remains the gold standard.

Further Statistical Explorations in R

Mastering the use of the polychoric coefficient is a significant step in accurate statistical modeling involving categorical data. It ensures that the measurement scale limitations do not distort the estimation of genuine relationships between variables. Researchers in fields requiring robust correlation estimates for survey data should integrate this technique into their analytical toolkit.

The principles explored here regarding latent variables and MLE extend far beyond simple correlation and are central to advanced techniques such as Factor Analysis and Item Response Theory. Understanding the limitations of standard correlations and knowing when to apply specialized measures like PCC is a hallmark of rigorous quantitative analysis.

The following tutorials explain how to calculate other common correlation coefficients in R: