

# Spearman's rank correlation coefficient

In statistics, **Spearman's rank correlation coefficient** or **Spearman's rho**, named after Charles Spearman and often denoted by the Greek letter ***ρ*** (rho) or as ***r<sub>s</sub>***, is a nonparametric measure of rank correlation (statistical dependence between the rankings of two variables). It assesses how well the relationship between two variables can be described using a monotonic function.

The Spearman correlation between two variables is equal to the Pearson correlation between the rank values of those two variables; while Pearson's correlation assesses linear relationships, Spearman's correlation assesses monotonic relationships (whether linear or not). If there are no repeated data values, a perfect Spearman correlation of +1 or −1 occurs when each of the variables is a perfect monotone function of the other.

Intuitively, the Spearman correlation between two variables will be high when observations have a similar (or identical for a correlation of 1) rank (i.e. relative position label of the observations within the variable: 1st, 2nd, 3rd, etc.) between the two variables, and low when observations have a dissimilar (or fully opposed for a correlation of −1) rank between the two variables.

Spearman's coefficient is appropriate for both continuous and discrete ordinal variables.<sup>[1][2]</sup> Both Spearman's ***ρ*** and Kendall's ***τ*** can be formulated as special cases of a more general correlation coefficient.

## Contents

**Definition and calculation**

**Related quantities**

**Interpretation**

**Example**

**Determining significance**

**Correspondence analysis based on Spearman's rho**

**Software Implementations**

**See also**

**References**

**Further reading**

**External links**

## Definition and calculation

The Spearman correlation coefficient is defined as the Pearson correlation coefficient between the rank variables.<sup>[3]</sup>

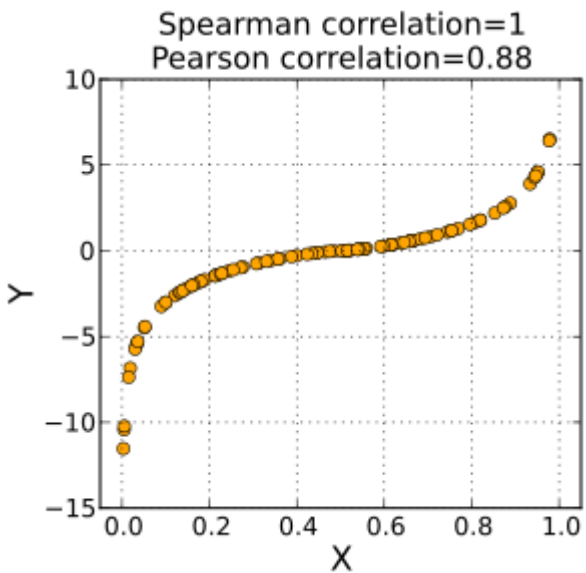
For a sample of size *n*, the *n* raw scores ***X<sub>i</sub>***, ***Y<sub>i</sub>*** are converted to ranks **rg** ***X<sub>i</sub>***, **rg** ***Y<sub>i</sub>***, and ***r<sub>s</sub>*** is computed from:

$$r_s = \rho_{\mathbf{rg}_X, \mathbf{rg}_Y} = \frac{\mathbf{cov}(\mathbf{rg}_X, \mathbf{rg}_Y)}{\sigma_{\mathbf{rg}_X} \sigma_{\mathbf{rg}_Y}}$$

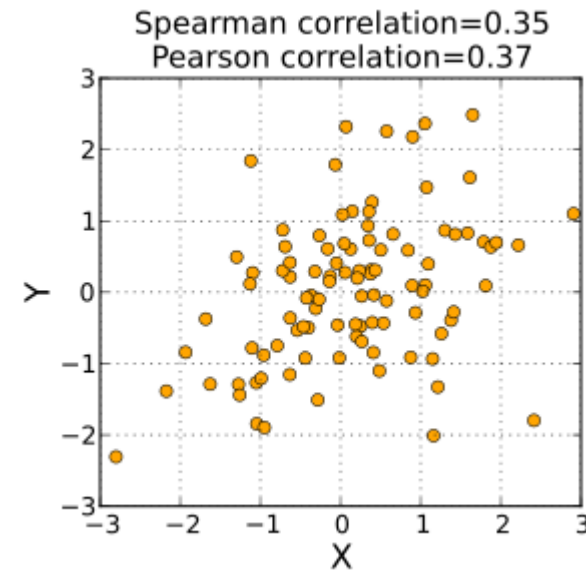
where

- ρ*** denotes the usual Pearson correlation coefficient, but applied to the rank variables.
- cov**(**rg** ***X***, **rg** ***Y***) is the covariance of the rank variables.
- σ***<sub>**rg** ***X***</sub> and ***σ***<sub>**rg** ***Y***</sub> are the standard deviations of the rank variables.

Only if all *n* ranks are *distinct integers*, it can be computed using the popular formula



A Spearman correlation of 1 results when the two variables being compared are monotonically related, even if their relationship is not linear. This means that all data-points with greater x-values than that of a given data-point will have greater y-values as well. In contrast, this does not give a perfect Pearson correlation.



When the data are roughly elliptically distributed and there are no prominent outliers, the Spearman correlation and Pearson correlation give similar values.

$$r_s = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}.$$

where

- d<sub>i</sub>* = **rg**(*X<sub>i</sub>*) − **rg**(*Y<sub>i</sub>*), is the difference between the two ranks of each observation.
- n* is the number of observations

Identical values are usually<sup>[4]</sup> each assigned fractional ranks equal to the average of their positions in the ascending order of the values, which is equivalent to averaging over all possible permutations.

If ties are present in the data set, the simplified formula above yields incorrect results: Only if in both variables all ranks are distinct, then ***σ<sub>rg<sub>X</sub></sub>σ<sub>rg<sub>Y</sub></sub>*** = **Var** ***rg<sub>X</sub>*** = **Var** ***rg<sub>Y</sub>*** = **(*n*<sup>2</sup> − 1)/12** (Calculated according to biased variance.). The first equation — normalizing by the standard deviation — may be used even when ranks are normalized to [0, 1] ("relative ranks") because it is insensitive both to translation and linear scaling.

The simplified method should also not be used in cases where the data set is truncated; that is, when the Spearman correlation coefficient is desired for the top X records (whether by pre-change rank or post-change rank, or both), the user should use the Pearson correlation coefficient formula given above.<sup>[5]</sup>

The standard error of the coefficient (*σ*) was determined by Pearson in 1907 and Gosset in 1920. It is

$$\sigma_{r_s} = \frac{0.6325}{\sqrt{n-1}}$$

## Related quantities

There are several other numerical measures that quantify the extent of statistical dependence between pairs of observations. The most common of these is the Pearson product-moment correlation coefficient, which is a similar correlation method to Spearman's rank, that measures the “linear” relationships between the raw numbers rather than between their ranks.

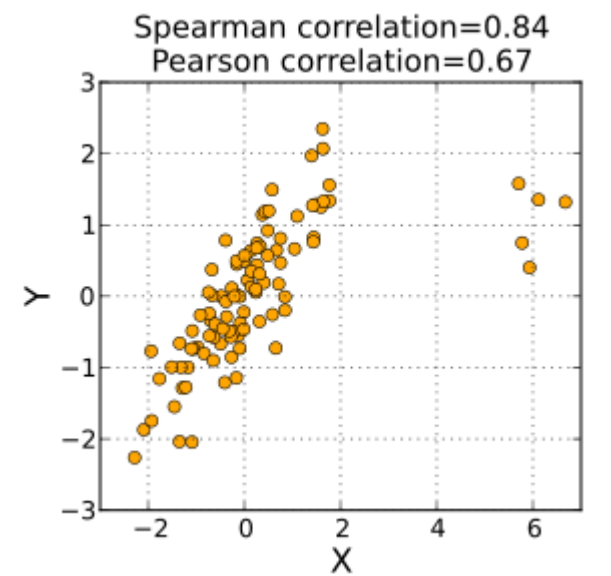
An alternative name for the Spearman rank correlation is the “grade correlation”;<sup>[6]</sup> in this, the “rank” of an observation is replaced by the “grade”. In continuous distributions, the grade of an observation is, by convention, always one half less than the rank, and hence the grade and rank correlations are the same in this case. More generally, the “grade” of an observation is proportional to an estimate of the fraction of a population less than a given value, with the half-observation adjustment at observed values. Thus this corresponds to one possible treatment of tied ranks. While unusual, the term “grade correlation” is still in use.<sup>[7]</sup>

## Interpretation

The sign of the Spearman correlation indicates the direction of association between *X* (the independent variable) and *Y* (the dependent variable). If *Y* tends to increase when *X* increases, the Spearman correlation coefficient is positive. If *Y* tends to decrease when *X* increases, the Spearman correlation coefficient is negative. A Spearman correlation of zero indicates that there is no tendency for *Y* to either increase or decrease when *X* increases. The Spearman correlation increases in magnitude as *X* and *Y* become closer to being perfect monotone functions of each other. When *X* and *Y* are perfectly monotonically related, the Spearman correlation coefficient becomes 1. A perfect monotone increasing relationship implies that for any two pairs of data values *X<sub>i</sub>*, *Y<sub>i</sub>* and *X<sub>j</sub>*, *Y<sub>j</sub>*, that *X<sub>i</sub>* − *X<sub>j</sub>* and *Y<sub>i</sub>* − *Y<sub>j</sub>* always have the same sign. A perfect monotone decreasing relationship implies that these differences always have opposite signs.

The Spearman correlation coefficient is often described as being "nonparametric". This can have two meanings. First, a perfect Spearman correlation results when *X* and *Y* are related by any monotonic function. Contrast this with the Pearson correlation, which only gives a perfect value when *X* and *Y* are related by a *linear* function. The other sense in which the Spearman correlation is nonparametric is that its exact sampling distribution can be obtained without requiring knowledge (*i.e.*, knowing the parameters) of the joint probability distribution of *X* and *Y*.

## Example



The Spearman correlation is less sensitive than the Pearson correlation to strong outliers that are in the tails of both samples. That is because Spearman's rho limits the outlier to the value of its rank.

#### Positive and negative Spearman rank correlations

In this example, the raw data in the table below is used to calculate the correlation between the IQ of a person with the number of hours spent in front of TV per week.

<b>IQ, <math>X_i</math></b>	<b>Hours of TV per week, <math>Y_i</math></b>
106	7
86	0
100	27
101	50
99	28
103	29
97	20
113	12
112	6
110	17

Firstly, evaluate  $d_i^2$ . To do so use the following steps, reflected in the table below.

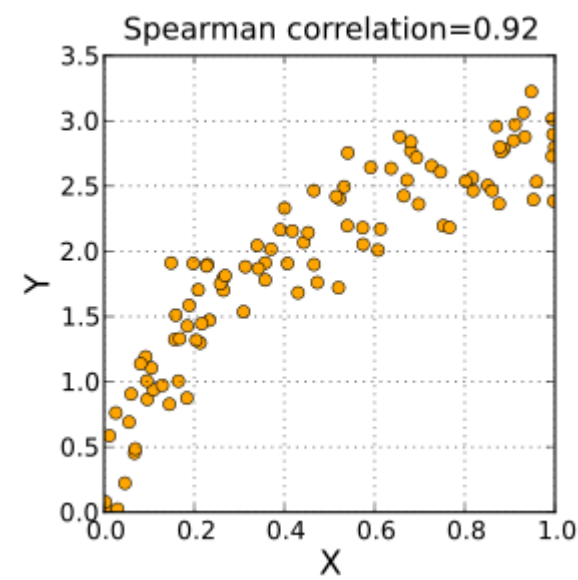
1. Sort the data by the first column ( $X_i$ ). Create a new column  $x_i$  and assign it the ranked values 1,2,3,... $n$ .
2. Next, sort the data by the second column ( $Y_i$ ). Create a fourth column  $y_i$  and similarly assign it the ranked values 1,2,3,... $n$ .
3. Create a fifth column  $d_i$  to hold the differences between the two rank columns ( $x_i$  and  $y_i$ ).
4. Create one final column  $d_i^2$  to hold the value of column  $d_i$  squared.

<b>IQ, <math>X_i</math></b>	<b>Hours of TV per week, <math>Y_i</math></b>	<b>rank <math>x_i</math></b>	<b>rank <math>y_i</math></b>	<b><math>d_i</math></b>	<b><math>d_i^2</math></b>
86	0	1	1	0	0
97	20	2	6	-4	16
99	28	3	8	-5	25
100	27	4	7	-3	9
101	50	5	10	-5	25
103	29	6	9	-3	9
106	7	7	3	4	16
110	17	8	5	3	9
112	6	9	2	7	49
113	12	10	4	6	36

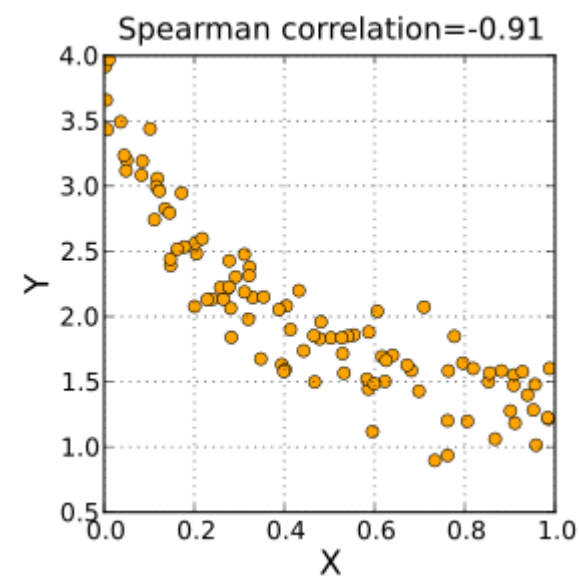
With  $d_i^2$  found, add them to find  $\sum d_i^2 = 194$ . The value of  $n$  is 10. These values can now be substituted back into the equation:  $\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$ , to give

$$\rho = 1 - \frac{6 \times 194}{10(10^2 - 1)}$$

which evaluates to  $\rho = -29/165 = -0.175757575...$  with a P-value = 0.627188 (using the t distribution).



A positive Spearman correlation coefficient corresponds to an increasing monotonic trend between X and Y.



A negative Spearman correlation coefficient corresponds to a decreasing monotonic trend between X and Y.

That the value is close to zero shows that the correlation between IQ and hours spent watching TV is very low, although the negative value suggests that the longer the time spent watching television the lower the IQ. In the case of ties in the original values, this formula should not be used; instead, the Pearson correlation coefficient should be calculated on the ranks (where ties are given ranks, as described above).

## Determining significance

One approach to test whether an observed value of *ρ* is significantly different from zero (*r* will always maintain −1 ≤ *r* ≤ 1) is to calculate the probability that it would be greater than or equal to the observed *r*, given the null hypothesis, by using a permutation test. An advantage of this approach is that it automatically takes into account the number of tied data values there are in the sample, and the way they are treated in computing the rank correlation.

Another approach parallels the use of the Fisher transformation in the case of the Pearson product-moment correlation coefficient. That is, confidence intervals and hypothesis tests relating to the population value *ρ* can be carried out using the Fisher transformation:

$$F(r) = \frac{1}{2} \ln \frac{1+r}{1-r} = \mathbf{arctanh}(r).$$

If *F*(*r*) is the Fisher transformation of *r*, the sample Spearman rank correlation coefficient, and *n* is the sample size, then

$$z = \sqrt{\frac{n-3}{1.06}} F(r)$$

is a z-score for *r* which approximately follows a standard normal distribution under the null hypothesis of statistical independence (*ρ* = 0).<sup>[8]</sup><sup>[9]</sup>

One can also test for significance using

$$t = r \sqrt{\frac{n-2}{1-r^2}}$$

which is distributed approximately as Student's t distribution with *n* − 2 degrees of freedom under the null hypothesis.<sup>[10]</sup> A justification for this result relies on a permutation argument.<sup>[11]</sup>

A generalization of the Spearman coefficient is useful in the situation where there are three or more conditions, a number of subjects are all observed in each of them, and it is predicted that the observations will have a particular order. For example, a number of subjects might each be given three trials at the same task, and it is predicted that performance will improve from trial to trial. A test of the significance of the trend between conditions in this situation was developed by E. B. Page<sup>[12]</sup> and is usually referred to as Page's trend test for ordered alternatives.

## Correspondence analysis based on Spearman's rho

Classic correspondence analysis is a statistical method that gives a score to every value of two nominal variables. In this way the Pearson correlation coefficient between them is maximized.

There exists an equivalent of this method, called grade correspondence analysis, which maximizes Spearman's rho or Kendall's tau.<sup>[13]</sup>

## Software Implementations

- R's statistics base-package implements the test  `cor.test(x, y, method = "spearman")` (http://stat.ethz.ch/R-manual/R-patched/library/stats/html/cor.test.html) in its "stats" package (also  `cor(x, y, method = "spearman")` will work, but without returning the p-value).

## See also

- Kendall tau rank correlation coefficient
- Chebyshev's sum inequality, rearrangement inequality (These two articles may shed light on the mathematical properties of Spearman's ρ.)
- Distance correlation

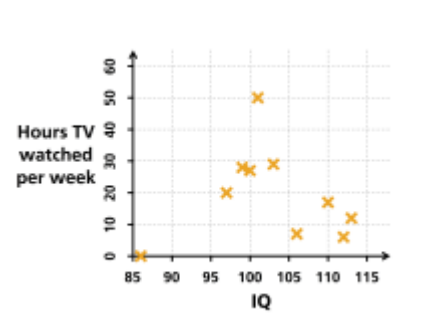


Chart of the data presented. It can be seen that there might be a negative correlation, but that the relationship does not appear definitive.