

Using Pearson's Correlation Coefficient

I. Correlation Coefficients

- Correlation coefficients are a class of statistical tests that provide standardized estimates of the strength of the relationship between 2 variables.
- The type of data (nominal, ordinal, interval/ratio) determines the statistic that you use.
 - For example Cramer's V (a.k.a, phi) is used for discrete data
- In general they represent the ratio of the variance shared between two variables (Covariance) divided by the total variability in both measures.
- Squaring a correlation coefficient tells us the % of the variance in the Dependent Variable that is explained by / predicted by / accounted for by variability in the Independent Variable

II. Pearson's r

- Used with two continuous variables
- Estimates the Strength and Direction of the Relationship between the IV and DV
- Scores range from -1 to 1
 - the sign of the coefficient indicates the direction of the relationship
 - 1 = perfect negative correlation (inverse relationship)
 - As one variable increases the other variable decreases.
 - 1 = perfect positive correlation (direct relationship)
 - As one variable increases the other variable increases.
 - the closer the absolute value of r is to 1 (regardless of whether it is positive or negative), the stronger the relationship between the two variables
 - the closer the value of r is to 0, the weaker the relationship is between the two variables.

Statistical Hypotheses: $H_0 : r = 0$
 $H_a : r \neq 0$

- we are testing whether or not the strength of the relationship between the IV and the DV differs from what we would expect by chance alone. That is, does it significantly differ from 0.

Example Research Question: What is the strength and direction of the association between the number of Vienna Sausages that a person eats in 24 hr period (measured in grams, range = 0-900) and the amount of time they spend doing impersonations of south park characters in that same 24 hr. period (measured in minutes, range = 0-1440).

III. Calculating Pearson's r .

Pearson's r : Allows us to test the strength of the association between two continuous variables. It represents a ratio of the Covariance (variance shared by two variables) and the total variance (covariance + unique variance). **df = n-2**

$$r = \frac{\Sigma XY - \frac{(\Sigma X)(\Sigma Y)}{n}}{\sqrt{\Sigma X^2 - \frac{(\Sigma X)^2}{n}} \sqrt{\Sigma Y^2 - \frac{(\Sigma Y)^2}{n}}} \quad \text{or} \quad r = \frac{\Sigma (X - \bar{X})(Y - \bar{Y})}{\sqrt{\Sigma (X - \bar{X})^2} \sqrt{\Sigma (Y - \bar{Y})^2}}$$

r^2 : The Coefficient of Determination : r represents a ratio of covariance to total variance, but if we want to know what percent of the variance in the DV is explained by variance in the IV, then we can square r to find out.

$1-r^2$ = the amount of variance (as a percentage) that is not accounted for. How much residual (left over)

variance there is. How much variance is due to error.

Error = The influence on the dependent variable that is attributed to sources other than the independent variable. Also called residual variance and unexplained variance.

IV. How Does r Work.

-What the r statistic does is divide up the variance or the deviation in the data set. The deviation of each observation (subject score) is made up of two parts.

Total Variance = (Variance Shared by Variables) + (Variance due to Error)
Covariance

Covariance = the amount of variance in (Y) that is explained (accounted for) by the variance in (X).

$$\text{Covariance}_{xy} = \frac{\sum (X - \bar{X})(Y - \bar{Y})}{n-1} \quad \text{Or} \quad \frac{\sum XY - [(\sum X)(\sum Y)/n]}{n - 1}$$

Variance due to Error = the amount of variance in (Y) that is not explained (shared with / accounted for) by the variance in (X).

What the r formula does is determine the total amount of Covariance and then divide this value by the total amount of variance.

-Pearson's r Formula

$$r = \frac{\sum XY - \frac{(\sum X)(\sum Y)}{n}}{\sqrt{\sum X^2 - \frac{(\sum X)^2}{n}} \sqrt{\sum Y^2 - \frac{(\sum Y)^2}{n}}} \quad r = \frac{\text{Covariance}}{\text{Total Variance for Each Variable}}$$

V. Significance of r .

- Compare the absolute value of the r - obtained, with the r - critical, for $df = n - 2$, at the $p = .05$ level.

- The table below presents the critical values of Pearson's r .

- If r - obtained $>$ r - critical_(df = n-2, p = .05, two tailed), then reject H_0 and Fail to reject H_A

- If r - obtained $<$ r - critical_(df = n-2, p = .05, two tailed), then reject H_A and Fail to reject H_0

- The table below presents One and Two tailed tests. For our purposes, we are really only interested in a Two Tailed test.

- There are occasions where a One Tailed test is permissible, but they are relatively rare and we don't need to worry about them at this point.

Table of Critical Values for Pearson's r .

<i>df</i>	Level of Significance for a One-Tailed Test					
	.10	.05	.025	.01	.005	.0005
	Level of Significance for a Two-Tailed Test					
	.20	.10	.05	.02	.01	.001
1	0.951	0.988	0.997	0.9995	0.9999	0.999999
2	0.800	0.900	0.950	0.980	0.990	0.999
3	0.687	0.805	0.878	0.934	0.959	0.991
4	0.608	0.729	0.811	0.882	0.917	0.974
5	0.551	0.669	0.755	0.833	0.875	0.951
6	0.507	0.621	0.707	0.789	0.834	0.925
7	0.472	0.582	0.666	0.750	0.798	0.898
8	0.443	0.549	0.632	0.715	0.765	0.872
9	0.419	0.521	0.602	0.685	0.735	0.847
10	0.398	0.497	0.576	0.658	0.708	0.823
11	0.380	0.476	0.553	0.634	0.684	0.801
12	0.365	0.457	0.532	0.612	0.661	0.780
13	0.351	0.441	0.514	0.592	0.641	0.760
14	0.338	0.426	0.497	0.574	0.623	0.742
15	0.327	0.412	0.482	0.558	0.606	0.725
16	0.317	0.400	0.468	0.542	0.590	0.708
17	0.308	0.389	0.456	0.529	0.575	0.693
18	0.299	0.378	0.444	0.515	0.561	0.679
19	0.291	0.369	0.433	0.503	0.549	0.665
20	0.284	0.360	0.423	0.492	0.537	0.652
21	0.277	0.352	0.413	0.482	0.526	0.640
22	0.271	0.344	0.404	0.472	0.515	0.629
23	0.265	0.337	0.396	0.462	0.505	0.618
24	0.260	0.330	0.388	0.453	0.496	0.607
25	0.255	0.323	0.381	0.445	0.487	0.597
26	0.250	0.317	0.374	0.437	0.479	0.588
27	0.245	0.311	0.367	0.430	0.471	0.579
28	0.241	0.306	0.361	0.423	0.463	0.570
29	0.237	0.301	0.355	0.416	0.456	0.562
30	0.233	0.296	0.349	0.409	0.449	0.554
40	0.202	0.257	0.304	0.358	0.393	0.490
60	0.165	0.211	0.250	0.295	0.325	0.408
120	0.117	0.150	0.178	0.210	0.232	0.294
∞	0.057	0.073	0.087	0.103	0.114	0.146

Adapted from Appendix 2 (Critical Values of t) using the square root of $[t^2/(t^2 + df)]$

Note: Critical Values for infinite degrees of freedom actually calculated at $df = 500$.