

## **A SIMPLER METHOD FOR TEACHING THE MEANING OF THE CORRELATION COEFFICIENT**

**by John N. Zorich, Jr.**

For over 100 years, the concept of the correlation coefficient (CC) has been taught to beginning students of statistical science but without serious reference to the equation on which it is based. Instead, the meaning of the CC has been explained using wordy generalities and textbook scatter plots — the CC being larger the less scattered the plot looks. Unfortunately, such generalities result in most students internalizing subtle misconceptions. Figs. 1, 2, and 3 demonstrate some of the difficulties that cannot be explained using classic teaching methods.

In Fig. 1, each of the four Data Sets (labeled A, B, C, and D) has a least squares linear regression (LSLR) straight line drawn through the raw data points. Each Data Set has 2 different Y values for each even X whole number from 2 through 18 (in Set D, the 2 different Y values at each X value are so close together that they appear as a single dot). In spite of the obvious differences between these data sets, they all yield the same high CC.

In Fig. 2, Data Sets B and C are, in effect, subsets of Set A. Set B is composed of the first six data points from Set A after subtracting 3.0 from each Y value. Likewise, Set C is the first three data points from Set A after subtracting 6.0. Notice that all three regression lines have the identical slope and have data points that lie at exactly the same distance from their regression line. It seems that a coefficient that purports to indicate correlation should indicate that these three data sets are, correlatively speaking, the same. But, as indicated in the figure, the larger the data set, the larger the CC.

In Fig. 3, we seem to have a contradiction to the conclusion reached regarding Fig. 2; that is, in Fig. 3, the more data points, the lower the CC, despite the fact that all three regression lines have the identical slope and have data points that lie at exactly the same distance from their regression line (as in Fig. 2, Sets B and C are subsets of Set A, offset by a value of 3.0 or 6.0, respectively).

The separation of CC meaning from the CC equation stems historically from the fact that the equations that appeared most often in textbooks were difficult to teach or understand. The first equations were developed in the late 1800s<sup>1</sup>; the most commonly cited one has been some version of the following:

$$CC = \frac{\sum (X - \bar{X})(Y - \bar{Y})}{\sqrt{\sum (X - \bar{X})^2 \sum (Y - \bar{Y})^2}} \quad \text{Equation 1} \quad (\text{the } \underline{\text{traditional equation}})$$

The “Short Method”<sup>2</sup> (Equation 2) was popularized in the early 1900s, and became more common after it was revised slightly and renamed the “computational form”<sup>3</sup> for use with sophisticated mechanical calculators and simple electronic ones:

## A Simpler Method for Teaching the Meaning of the Correlation Coefficient

Copyright 2012, by John N. Zorich Jr., Zorich Consulting & Training, www.johnzorich.com

---

$$CC = \frac{N \sum XY - \sum X \sum Y}{\left( \sqrt{N \sum X^2 - (\sum X)^2} \right) \left( \sqrt{N \sum Y^2 - (\sum Y)^2} \right)} \quad \text{Equation 2}$$

A simpler equation<sup>4</sup> (Equation 3) was developed, but it didn't have much pedagogic value. In this equation,  $S_x$  and  $S_y$  are the standard deviation of the X and Y data, respectively, and "slope" is the slope of the linear regression line calculated by the method of least squares:

$$CC = \frac{(\text{slope})S_x}{S_y} \quad \text{Equation 3}$$

The ratio of slope to  $S_y$  in this equation is helpful in explaining the *lack* of dependence of the CC on slope, since a large CC can result from either a large slope, a large  $S_x$ , and/or a small  $S_y$ . This equation is unable to explain what the CC is, but is wonderful for explaining what the CC is not.

About 1985, I thought I'd developed a new, more instructive equation. Alas, a few years later, I found my "new" equation at the end of an appendix to a 1961 introductory statistics book written by someone else.<sup>5</sup>

I discovered my "new" equation within the equation for the Coefficient of Determination (CD). In the CD equation (see next),  $Y_e$  represents the Y values calculated for each X value by using the least squares linear regression analysis equation ( $Y_e = a + bX$ ), and  $Y_i$  represents the raw Y data. One  $Y_e$  value is calculated for each  $Y_i$  value. As always in least squares linear regression, the mean of the  $Y_e$  data is the same as that of the  $Y_i$ , and so it is not subscripted in the equation below:

$$CD = CC^2 = \frac{\sum (Y_e - \bar{Y})^2}{\sum (Y_i - \bar{Y})^2} \quad \text{Equation 4}$$

Dividing top and bottom of the fraction by  $N-1$  (where  $N$  is the number of  $Y_i$  data points), I discovered an equation that is the ratio of two sample variances:

$$CC^2 = \frac{\text{Variance}(Y_e)}{\text{Variance}(Y_i)} \quad \text{Equation 5}$$

After taking the square root of both sides, I found an equation containing the absolute value of the CC on one side, and the ratio of two sample standard deviations on the other:

$$|CC| = \frac{\text{StdDeviation}(Y_e)}{\text{StdDeviation}(Y_i)} \quad \text{Equation 6 (the "new" equation)}$$

## A Simpler Method for Teaching the Meaning of the Correlation Coefficient

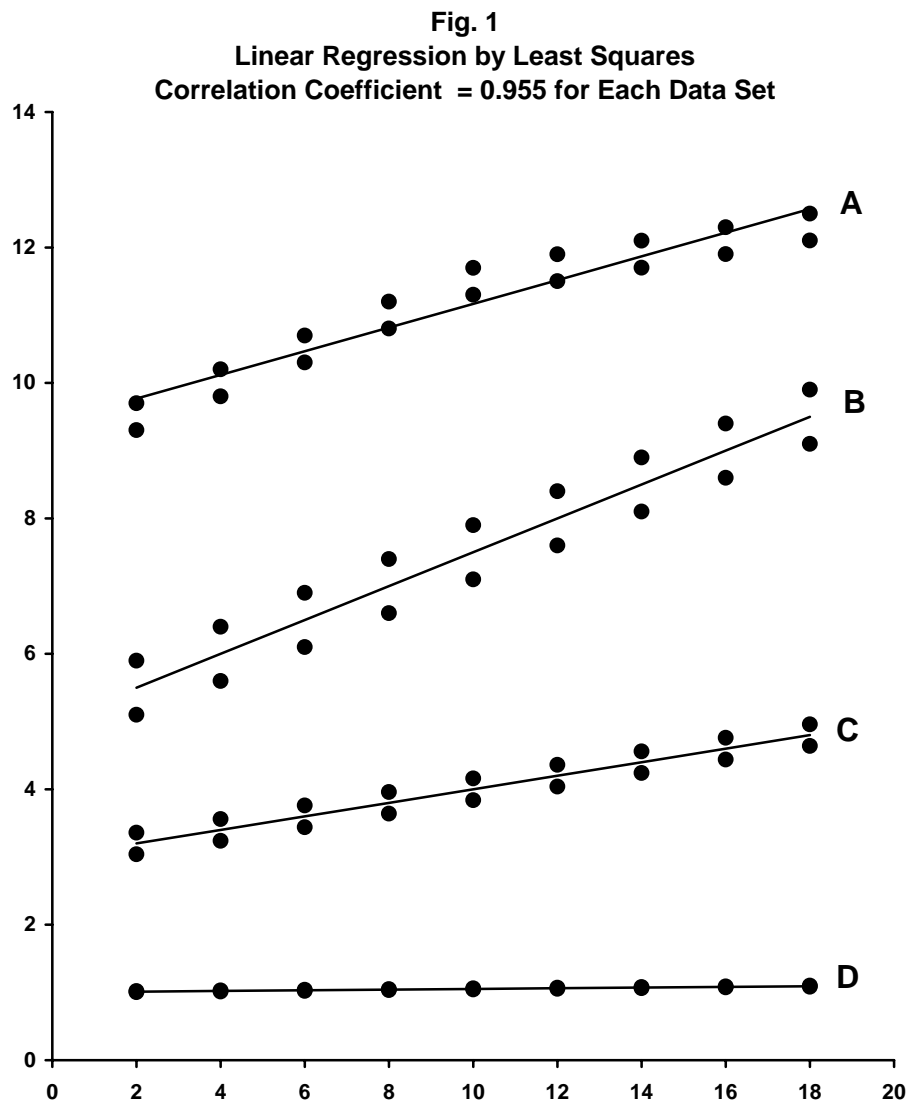
Copyright 2012, by John N. Zorich Jr., Zorich Consulting & Training, [www.johnzorich.com](http://www.johnzorich.com)

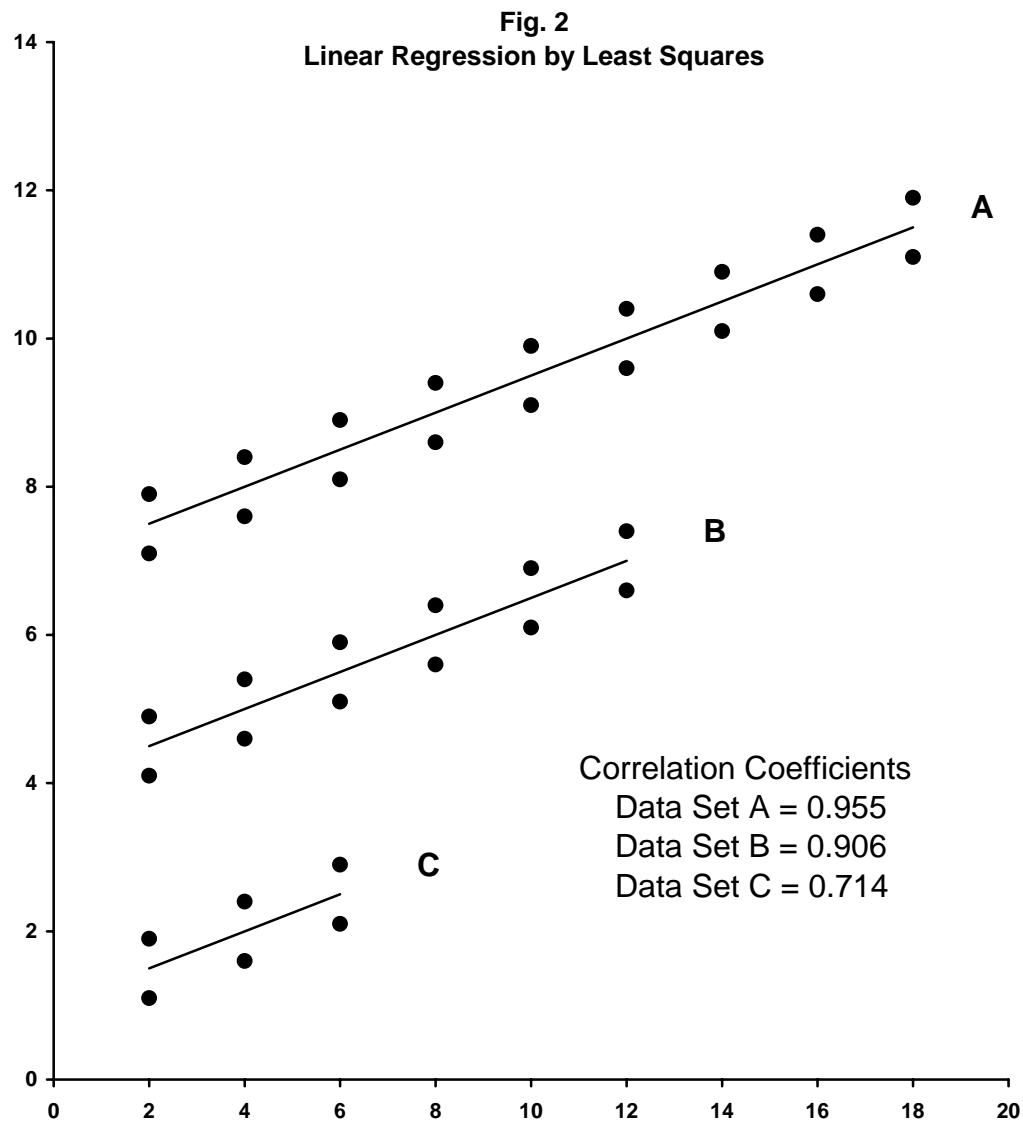
---

Although this equation does not calculate the sign of the CC, this is not a limitation. In LSLR, the sign of the CC is always the same as that of the regression coefficient (“b” in the LSLR equation  $Y = a + bX$ ) — that is, if the slope is negative, so is the CC, and vice versa, as easily seen from Equation 3, above. Thus, the sign of the CC has no meaning independent of the regression slope, and so the only unique aspect of the CC is its absolute value, which the “new” equation calculates.

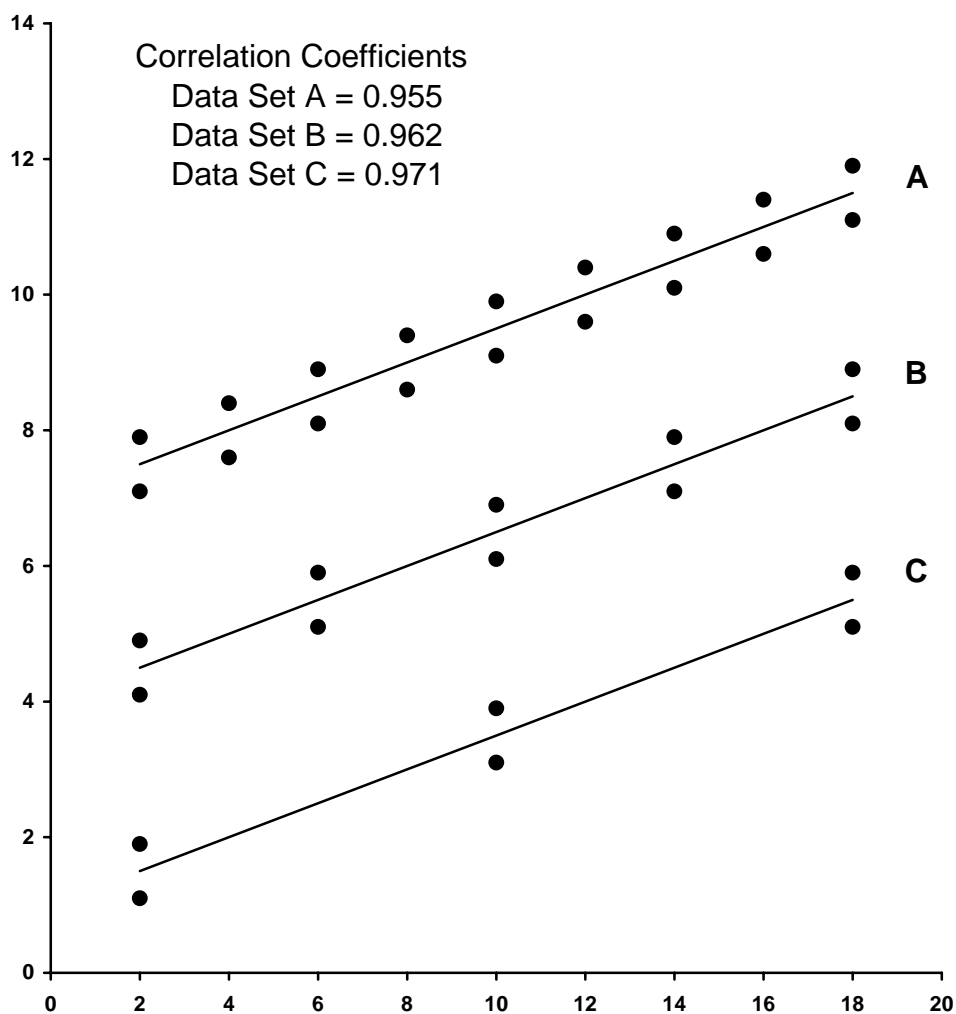
Calculation of the CC using the new equation is shown by example in Table 1. Not only can this new equation be used to easily explain the surprising CC results shown in Figs. 1, 2, and 3, but it can also be used to explain other interesting facts, such as:

1. The correlation coefficient can never equal exactly 1.000, unless all the  $Y_i$ 's form a perfectly straight line — which is the only case in which the standard deviations of  $Y_e$  and  $Y_i$  are identical.
2. The CC can never equal exactly 0.000, unless the standard deviation of  $Y_e$  is also zero — which would occur only if the calculated linear regression line were perfectly horizontal.
3. The CC represents the fraction of the total variation in  $Y_i$ , as measured in units of standard deviation, that can be explained by a linear relationship between  $Y_i$  and  $X$ . The larger the CC, the larger the fraction of the  $Y_i$  variation which can be explained this way. The remaining variation can't be explained, at least not by the CC.





**Fig. 3**  
**Linear Regression by Least Squares**



# A Simpler Method for Teaching the Meaning of the Correlation Coefficient

Copyright 2012, by John N. Zorich Jr., Zorich Consulting & Training, [www.johnzorich.com](http://www.johnzorich.com)

<b>Table 1.</b> <b>Calculation of the Correlation Coefficient using the “New Equation”</b> (Ye is calculated using the regression equation at the bottom of this table)		
<b>X raw data</b>	<b>Yi raw data</b>	<b>Ye, calculated</b>
6	10	9.7609
7	11	10.2065
7	10	11.5435
8	11	11.5435
9	12	11.5435
10	12	11.9891
10	12	12.4348
12	13	12.8804
13	14	13.3261
15	14	13.7717
	Standard Deviation = 1.4491	Standard Deviation = 1.4023
Least Squares Linear Regression Equation is $Y_e = 7.2221 + 0.4823X$ Correlation Coefficient (CC) using the “Traditional Equation” = <b>0.9677</b> CC using the “New Equation,” $(\text{StdDev } Y_e) / (\text{StdDev } Y_i) =$ <b>0.9677</b>		

<sup>1</sup> S. M. Stigler, *The History of Statistics*, 1986 (Belknap Press, Cambridge MA), chapter 9.

<sup>2</sup> J. G. Smith, *Elementary Statistics*, 1934 (Henry Holt & Co., New York) p. 374.

<sup>3</sup> H. L. Alder, E. B. Roessler, *Introduction to Probability and Statistics*, 6th ed., 1977 (W. H. Freeman & Co., San Francisco) p. 230.

<sup>4</sup> *Ibid*, p. 231. Alder & Roessler use different symbols than are used here.

<sup>5</sup> W. J. Reichmann, *Use and Abuse of Statistics*, 1961 (Oxford University Press, New York) p. 306.