

A Geometrical Meaning for the Correlation Coefficient

Murray L Lauber

Having the privilege of teaching many mathematics strands has allowed me to make connections between concepts that on the surface seem unrelated. High school and undergraduate math students have this opportunity but may not have the time or the intimate knowledge of the subject matter to form such connections. A reward of repeatedly teaching the same math courses is that one's knowledge of the subject matter deepens to the point where, with some exploration, such connections can become apparent. Additional rewards are modelling such exploration for students and encouraging them to explore on their own.

This article describes the fundamental connection between the concept of the correlation coefficient from statistics and that of the angle between two vectors from linear algebra. That connection became apparent to me over a few years while teaching vectors in linear algebra and, at the same time, some elementary statistics in a precalculus course. It initially sprouted from a concept that had incubated when I was a student in a statistics course many years earlier. In the chapter of the course textbook pertaining to the correlation coefficient, Ferguson (1981, 132) describes how the correlation coefficient is related to the angular separation between two regression lines. The ensuing discussion is general enough to leave room for questions and to invite exploration. In fact, Ferguson's observations seemed inaccurate because a full mathematical explanation was not given. At the least, they lodged in the back of my mind as a kind of healthy dissonance. They were not completely resolved until I taught a linear algebra course where the concept of the angle between two n -dimensional vectors was fully developed as the generalization of the geometrical angle between two 2- or 3-dimensional vectors. The angle between a pair of 2- or 3-dimensional vectors can be visualized intuitively and is easily calculated using simple trigonometry. The angle between two n -dimensional vectors is then defined as a generalization of the intuitive notions applying to 2- or 3-dimensional vectors.

What follows is the full development of a geometrical meaning for the correlation coefficient based on the notions of the previous paragraph. It is related to Ferguson's observations but, given an understanding of some basic concepts of vectors, seems more elegant in its simplicity.

The Correlation Coefficient— A Brief Review

The peripheral correlation coefficient is a precise comparison of two sets of scores that measures the degree to which corresponding scores deviate from their respective means. Do the sizes and directions of the deviations of corresponding data elements from their respective means tend to correspond? If so, the correlation coefficient will be high (close to 1). Does there appear to be little relationship between how corresponding data elements deviate from their respective means in the two sets of scores? If so, the correlation coefficient will be low (close to 0). Do the deviations from their respective means for corresponding elements tend to be in opposite directions (scores above the mean for the one data set correspond to scores below the mean for the other set, and vice versa)? If so, the correlation coefficient will be negative (perhaps as negative as -1). Consider the following simple example for two sets of scores, x and y .

x	y
1	2
2	4
3	6
4	8

[1]

Intuitively, these two sets of data are as closely related as any two distinct sets of data can be; therefore, the correlation coefficient should be 1. This will be demonstrated shortly.

The correlation coefficient may be defined as the ratio of the average of the sum of products of the

deviations of corresponding elements from their respective means to the product of the standard deviations of the two sets of scores. Consider the following two sets of scores¹.

$$x = (x_1, x_2, \dots, x_n), y = (y_1, y_2, \dots, y_n)$$

The correlation coefficient r between these two sets of scores is defined formulaically as follows:

$$r = \frac{(x_1 - \bar{x})(y_1 - \bar{y}) + (x_2 - \bar{x})(y_2 - \bar{y}) + \dots + (x_n - \bar{x})(y_n - \bar{y})}{ns_x s_y} \quad [2]$$

In this formulation:

- \bar{x} is the mean for the set x
- \bar{y} is the mean for the set y
- s_x is the standard deviation for the set x
- s_y is the standard deviation for the set y
- n is the number of scores in each data set

Definition [2] readily shows that the numerator will be large and positive if corresponding scores from x and y deviate proportionately in the same direction from their respective means. On the other hand, it will be large and negative if corresponding scores from x and y deviate proportionately in opposite directions from their respective means. And it will be small if there is little connection between how corresponding scores from x and y deviate from their respective means. The numerator alone, though, would not adequately define any measure of comparison between two sets of scores. We would be left with the questions, "How large is large?" and "How small is small?" But definition [2] taken altogether is ingenious in that dividing by $ns_x s_y$ ensures that the value of r is between -1 and 1 for any two sets of scores with 1 representing the highest possible positive correlation and -1 representing the lowest possible negative correlation. The proof is not included here but can be formed using the definitions of s_x , s_y , and r .

By way of illustration, Table A shows the calculations used in determining r for example [1].

Recall that we had already anticipated that the value of r for this case should be 1 . Note here that $\bar{x} = 2.5$, $\bar{y} = 5$ and $n = 4$. From Table A, we have

$$s_x = \sqrt{\frac{\sum (x - \bar{x})^2}{n}} = \frac{\sqrt{5}}{2}, s_y = \sqrt{\frac{\sum (y - \bar{y})^2}{n}} = \sqrt{5}$$

$$\text{Then } r = \frac{\sum (x - \bar{x})(y - \bar{y})}{ns_x s_y} = \frac{10}{4 \left(\frac{\sqrt{5}}{2} \right) (\sqrt{5})} = 1$$

This is as we expected.

Correlation Coefficients from Standard Scores

When comparing two data sets, it often helps to first convert the raw scores into standard scores or z -scores. The z -score of a particular score in a set of raw scores is the measure of how many standard deviations the raw score is above or below the mean. Suppose, for example, that for a set of scores x , the mean and standard deviation are $\bar{x} = 10$ and $s_x = 2$, respectively. Then a raw score of 12 would have a z -score of 1 because it is exactly one standard deviation above the mean. In general, the z -score, z_x , of a particular raw score x from the set of scores x where the mean is \bar{x} and the standard deviation is s_x is defined as

$$z_x = \frac{x - \bar{x}}{s_x}$$

The formulaic representation for the correlation coefficient r is simpler when standard scores are used. Recall that

$$r = \frac{\sum (x - \bar{x})(y - \bar{y})}{ns_x s_y} \quad [3]$$

Since $z_x = \frac{x - \bar{x}}{s_x}$ and $z_y = \frac{y - \bar{y}}{s_y}$, we have

$$r = \frac{\sum z_x z_y}{n} \quad [4]$$

Table A

x	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})^2$	$(y - \bar{y})^2$	$(x - \bar{x})(y - \bar{y})$
1	2	-1.5	-3	2.25	9	4.5
2	4	-0.5	-1	.25	1	.5
3	6	0.5	1	.25	1	.5
4	8	1.5	3	2.25	9	4.5
				$\sum (x - \bar{x})^2 = 5$	$\sum (y - \bar{y})^2 = 20$	$\sum (x - \bar{x})(y - \bar{y}) = 10$

This formula for r will be revisited after the concept of the angle between two vectors has been fully developed. It is most useful as a theoretical tool for developing other relationships. However, by way of illustration, it is applied to example [1] in Table B below. Recall that in this example, $\bar{x} = 2.5$, $\bar{y} = 5$, $s_x = \frac{\sqrt{5}}{2}$, and $s_y = \sqrt{5}$. By way of illustration, the values of z_x , z_y , and $z_x z_y$ in the first row were computed as follows.

$$z_x = \frac{-1.5}{\frac{\sqrt{5}}{2}} = \frac{-3}{\sqrt{5}}, z_y = \frac{-3}{\sqrt{5}}, \text{ and } z_x z_y = \left(\frac{-3}{\sqrt{5}}\right)\left(\frac{-3}{\sqrt{5}}\right) = \frac{9}{5}$$

Table B

x	y	$x - \bar{x}$	$y - \bar{y}$	z_x	z_y	$z_x z_y$
1	2	-1.5	-3	$-3\sqrt{5}$	$-3\sqrt{5}$	9/5
2	4	-0.5	-1	$-1\sqrt{5}$	$-1\sqrt{5}$	1/5
3	6	0.5	1	$1\sqrt{5}$	$1\sqrt{5}$	1/5
4	8	1.5	3	$3\sqrt{5}$	$3\sqrt{5}$	9/5
						$\sum z_x z_y = 4$

Using the results from the table, $r = \frac{\sum z_x z_y}{n} = \frac{4}{4} = 1$.

One other concept pertaining to z-scores will be needed to show the relationship between the correlation coefficient and the angle between two vectors. It is that of the magnitude of the vector formed by the z-scores of a data set. This notion will be easy to formulate but must await some basic concepts pertaining to vectors.

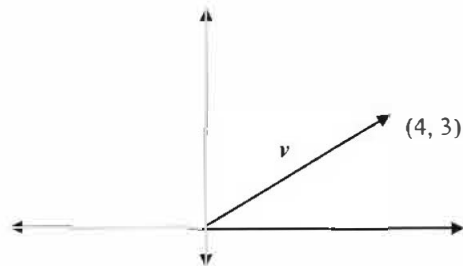
Vectors and Their Relevant Properties

What Is a Vector?

A vector is a directed line segment. A vector in the Cartesian plane is called a 2-dimensional geometric vector; a vector in Cartesian 3-space is called a 3-dimensional geometric vector. If the vector's initial point is at the origin of the Cartesian coordinate system, then the vector is in standard form. A vector not in standard form with initial point A and terminal point B is denoted \overline{AB} (in bold case). For convenience a vector may also be denoted as a single letter in bold case; for example, \mathbf{v} . If a 2- or 3-dimensional vector is in standard form, then it is determined by its terminal point. This leads to the following algebraic definitions for these vectors: a 2-dimensional vector is

an ordered pair of real numbers (a, b) ; a 3-dimensional vector is an ordered triple of real numbers (a, b, c) . Two- and 3-dimensional vectors can be represented geometrically. For example, the vector $\mathbf{v} = (4, 3)$ is illustrated in Figure 1.

Figure 1



Although we cannot picture more than three dimensions, the notions pertaining to algebraic vectors can be extended to any number of dimensions. An n -dimensional vector is defined as an ordered n -tuple (x_1, x_2, \dots, x_n) of real numbers. An n -dimensional vector is said to have n components. The i^{th} component is x_i .

The Length or Magnitude of a Vector

The length of a 2- or 3-dimensional vector can be determined easily using the formulas for the distance between a pair of points in 2- or 3-space, respectively. For example, the length of the vector $\mathbf{v} = (4, 3)$ in Figure 1 is $\|\mathbf{v}\| = \sqrt{4^2 + 3^2} = 5$. The magnitude of an algebraic vector is defined as being equal to the length of its corresponding geometric vector. So the terms *length* and *magnitude* are interchangeable. If $\mathbf{v} = (a, b)$, $a, b \in \mathbf{R}$, where \mathbf{R} is the set of real numbers, then the magnitude of \mathbf{v} , $\|\mathbf{v}\|$, is defined by $\|\mathbf{v}\| = \sqrt{a^2 + b^2}$; if $\mathbf{v} = (a, b, c)$, $a, b, c \in \mathbf{R}$, then $\|\mathbf{v}\| = \sqrt{a^2 + b^2 + c^2}$. These notions can be extended to n -dimensional vectors: if $\mathbf{v} = (x_1, x_2, \dots, x_n)$, then

$$\|\mathbf{v}\| = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2} \quad [5]$$

Of course, one cannot picture (at least in a sober state) the length of an n -dimensional vector if $n > 3$, but this definition is a reasonable abstraction consistent with our intuitive understanding of the lengths of 2- and 3-dimensional vectors.

The Inner (Dot Product) of Two Vectors

A number of operations are defined on vectors. Among them are two important products that involve pairs of vectors: the inner product and the cross product. Both have important applications as well as theoretical value. The one of relevance here is the inner product because it is useful in defining the

angle between two vectors. If $x = (x_1, x_2, \dots, x_n)$ and $y = (y_1, y_2, \dots, y_n)$, then the inner product $x \cdot y$ of x and y is

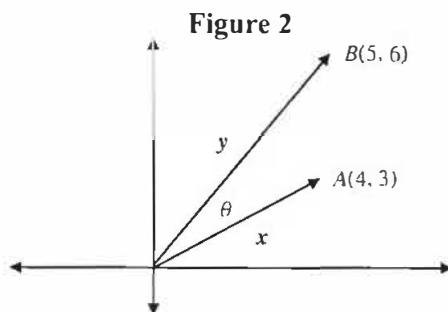
$$x \cdot y = x_1 y_1 + x_2 y_2 + \dots + x_n y_n = \sum x_i y_i \quad [6]$$

For example, if $x = (4, 3)$ and $y = (5, 6)$, then $x \cdot y = 4 \cdot 5 + 3 \cdot 6 = 38$.

The Angle Between Two Vectors

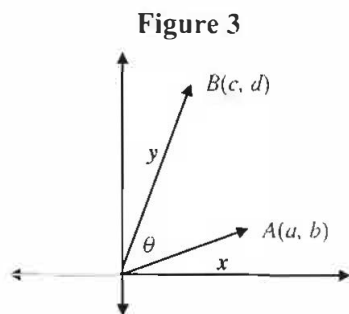
Consider the vectors $x = (4, 3)$ and $y = (5, 6)$ as illustrated in Figure 2. One can use the law of cosines to determine the angle θ between x and y :

$$\begin{aligned} \|AB\|^2 &= \|x\|^2 + \|y\|^2 - 2\|x\| \|y\| \cos \theta \\ \Rightarrow (5-4)^2 + (6-3)^2 &= (4^2 + 3^2) + (5^2 + 6^2) - \\ & 2 \sqrt{4^2 + 3^2} \sqrt{5^2 + 6^2} \cos \theta \\ \Rightarrow 10 &= 86 - 2 \cdot 5 \sqrt{61} \cos \theta \\ \Rightarrow \cos \theta &= \frac{76}{10\sqrt{61}} \\ \Rightarrow \theta &\cong 13.32^\circ \end{aligned}$$



Consider the general case for the angle between a pair of 2-dimensional vectors $x = (a, b)$ and $y = (c, d)$ in standard position as illustrated in Figure 3. Then, as in the previous example, $\|AB\|^2 = \|x\|^2 + \|y\|^2 - 2\|x\| \|y\| \cos \theta$

$$\begin{aligned} \Rightarrow \cos \theta &= (\|x\|^2 + \|y\|^2 - \|AB\|^2) / (2\|x\| \|y\|) \\ \Rightarrow \cos \theta &= (a^2 + b^2 + c^2 + d^2 - ((c-a)^2 + (d-b)^2)) / \\ & (2\|x\| \|y\|) \\ \Rightarrow \cos \theta &= (a^2 + b^2 + c^2 + d^2 - c^2 + 2ac - a^2 - d^2 + \\ & 2bd - b^2) / (2\|x\| \|y\|) \\ \Rightarrow \cos \theta &= (ac + bd) / (\|x\| \|y\|) \\ \Rightarrow \cos \theta &= (x \cdot y) / (\|x\| \|y\|) \quad [7] \end{aligned}$$



The result [7] provides a simple way of thinking about the angle θ between a pair of 2-dimensional vectors: the cosine of θ is just the inner product of the two vectors divided by the product of their magnitudes. Consider again the two vectors $x = (4, 3)$ and $y = (5, 6)$. Using [7] the angle θ between the two vectors is given by

$$\begin{aligned} \cos \theta &= (x \cdot y) / (\|x\| \|y\|) = (4 \cdot 5 + 3 \cdot 6) / (\sqrt{4^2 + 3^2} \sqrt{5^2 + 6^2}) \\ &= 38 / (5\sqrt{61}). \end{aligned}$$

This is the same value as that obtained earlier by more laborious methods.

The result [7] applies to 3-dimensional vectors as well. This can be seen by applying the law of cosines to a pair of 3-dimensional vectors $x = (a, b, c)$ and $y = (d, e, f)$. The steps are the same as those used above for 2-dimensional vectors. Verification of this result is left to the reader.

Although one cannot visualize the angle between two n -dimensional vectors for $n > 3$, it is reasonable to think of the angle between such a pair of vectors as a generalization of the angle between 2- or 3-dimensional vectors. This leads to the following definition. If $x = (x_1, x_2, \dots, x_n)$ and $y = (y_1, y_2, \dots, y_n)$ are any pair of n -dimensional vectors, then the angle θ between them is defined by

$$\cos \theta = (x \cdot y) / (\|x\| \|y\|) \quad [8]$$

Consider a simple example of a pair of 5-dimensional vectors $x = (1, 2, 3, 4, 5)$ and $y = (2, 4, 6, 8, 10)$. The vectors x and y have an obvious intuitive relationship to each other. In the precise language of vector algebra, y is said to be a scalar multiple of x . In general, a vector y is said to be a scalar multiple of vector x if each component of y is obtained from the corresponding component of x by multiplying by the same constant or scalar. In this case the constant is 2 and we write $y = 2x$. It is easy to appreciate why two n -dimensional vectors that are positive scalar multiples of each other are defined to have the same direction. Thus, in the above example, the vectors x and y should have the same direction and the angle between them should be 0° . Using definition [8] as follows yields a result that is consistent with this.

$$\begin{aligned} \cos \theta &= (x \cdot y) / (\|x\| \|y\|) \\ \Rightarrow \cos \theta &= (1 \cdot 2 + 2 \cdot 4 + 3 \cdot 6 + 4 \cdot 8 + 5 \cdot 10) / \\ & (\sqrt{1^2 + 2^2 + \dots + 5^2} \sqrt{2^2 + 4^2 + \dots + 10^2}) \\ \Rightarrow \cos \theta &= 110 / (\sqrt{55} \sqrt{220}) = 110 / (\sqrt{110^2}) = 1 \\ \Rightarrow \theta &= 0^\circ \end{aligned}$$

Data Sets as Vectors

With this framework, it is easy to see that a set of data can be represented as a vector. Consider the two

data sets $x = (x_1, x_2, \dots, x_n)$ and $y = (y_1, y_2, \dots, y_n)$ represented as vectors. Then it should be possible to formulate the correlation coefficient in terms of the vector concepts outlined in section 4 above. It turns out that the formulation is simpler if each set of scores is first converted to standard form; that is, each x_i and y_i is first converted to a z -score. We will refer to these vectors as the standard score vectors of x and y and denote them $z_x = (z_{x_1}, z_{x_2}, \dots, z_{x_n})$ and $z_y = (z_{y_1}, z_{y_2}, \dots, z_{y_n})$, respectively. Using [8] the angle θ between z_x and z_y is given by

$$\begin{aligned} \cos \theta &= (z_{x_1} z_{y_1} + z_{x_2} z_{y_2} + \dots + z_{x_n} z_{y_n}) / (\|z_x\| \|z_y\|) \\ \Rightarrow \cos \theta &= z_x \cdot z_y / (\|z_x\| \|z_y\|) \\ \Rightarrow \cos \theta &= \sum z_x z_y / (\|z_x\| \|z_y\|) \end{aligned} \quad [9]$$

We encountered the numerator of the right side of [9] earlier: it is also the numerator of the correlation coefficient in [4]. Let us examine the denominator $\|z_x\| \|z_y\|$. It can be shown that for the standard score vector z_x of any n -dimensional vector x , $\|z_x\| = \sqrt{n}$ as follows. Note that

$$\|z_x\| = \sqrt{z_{x_1}^2 + z_{x_2}^2 + \dots + z_{x_n}^2} = \sqrt{\sum z_{x_i}^2}$$

$$\text{But } z_{x_i} = \frac{x_i - \bar{x}}{s_i} \text{ and } s_i = \sqrt{\frac{\sum (x - \bar{x})^2}{n}}$$

$$\text{So } z_{x_i} = \frac{x_i - \bar{x}}{\sqrt{\frac{\sum (x - \bar{x})^2}{n}}} = \frac{(x_i - \bar{x})\sqrt{n}}{\sqrt{\sum (x - \bar{x})^2}}$$

$$\text{Then from [5], } \|z_x\| = \sqrt{\frac{(x_1 - \bar{x})^2 n}{\sum (x - \bar{x})^2} + \frac{(x_2 - \bar{x})^2 n}{\sum (x - \bar{x})^2} + \dots + \frac{(x_n - \bar{x})^2 n}{\sum (x - \bar{x})^2}}$$

$$\Rightarrow \|z_x\| = \sqrt{\frac{\sum (x - \bar{x})^2 n}{\sum (x - \bar{x})^2}} = \sqrt{n} \quad [10]$$

Since z_x and z_y are both standard score vectors, $\|z_x\| = \sqrt{n}$ and $\|z_y\| = \sqrt{n}$. Thus [9] becomes $\cos \theta = \frac{\sum z_{x_i} z_{y_i}}{n}$ [11]

The right side of [11] is the correlation coefficient between the set of scores x and y shown in formula [4]. Thus we have the result that has been the object of this article: the correlation coefficient between two sets of scores is just the cosine of the angle between their standard form vectors.

Applying formula [11] to the special cases where $\theta = 0^\circ, 90^\circ$ and 180° and noting that $\cos 0^\circ = 1$, $\cos 90^\circ = 0$ and $\cos 180^\circ = -1$ yields the following intriguing results about the value of the correlation coefficient r between the standard score vectors of two sets of scores:

- * $r = 1$ if and only if the standard score vectors have the same direction.
- * $r = -1$ if and only if the standard score vectors are in opposite directions.

* $r = 0$ if and only if the standard score vectors are perpendicular.

* r has a value between 0 and 1 if and only if the standard score vectors are somewhere between perpendicular and in the same direction.

* r has a value between 0 and -1 if and only if the standard score vectors are somewhere between perpendicular and in opposite directions.

Conclusion

This article demonstrates the relationship between correlation coefficient and the angle between two vectors. The beauty of this relationship is that it provides a simple geometrical meaning for the correlation coefficient that appeal to the intuition. There is also beauty and satisfaction in the processes underlying the discovery and development of this relationship. The exploratory and deductive methods used illustrate how mathematical connections are discovered and verified. Mathematics teachers who look for connections are in a good position to uncover such connections by virtue of the intimate knowledge of the subject matter that accompanies teaching. Further, they can model both the excitement and the discipline that is involved in carrying the discovery process to its conclusion. Teachers who are captivated by the exploration process will find ways to allow students to be captivated as well.

Notes

1. The two sets of data are presented here in vector notation; that is, as ordered n -tuples. This is a convenient notation and appropriate for the purposes of this article.

2. The formulations for standard deviation and the correlation coefficient used in this article are those pertaining to a whole population rather than a sample. Using n rather than the usual $n-1$ makes the demonstration of the relationship between the correlation coefficient and the angle between two vectors more transparent. But it is possible to demonstrate the relationship using $n-1$ as well.

3. The reader will recall that the standard deviation s_i for the data set $x = (x_1, x_2, \dots, x_n)$ is a measure of how the data is distributed about the mean \bar{x} . It is defined as follows.

$$s_i = \sqrt{\frac{\sum (x - \bar{x})^2}{n}}$$

and can be described as the root of the mean of the squares of the deviations of the individual scores from the mean.

Reference

Ferguson, G A. 1981. *Statistical Analysis in Psychology and Education*. 5th ed. New York: McGraw-Hill.

Note: The vector and statistics concepts underlying and related to this article can be found in any introductory linear algebra textbook or introductory statistics textbook.

Murray L Lauber is an associate professor of mathematics in the Augustana Faculty of the University of Alberta in Camrose, Alberta. He has taught a variety of courses including precalculus, introductory calculus, linear algebra, discrete mathematics, history of mathematics and higher arithmetic. He believes that mathematics is a potent tool for expanding the intellectual capacities of all

students. He shares his passion for mathematical exploration and discovery with his students and his colleagues and peers, particularly through his writing. Many of his articles, in which he shares intriguing mathematical relationships that he has uncovered during the course of his teaching, have been published in delta-K and in the Mathematics Teacher, a publication of the National Council of Teachers of Mathematics.