



THE GEOMETRIC INTERPRETATION OF THE CORRELATION RATIO*

Renato Leoni

Dipartimento Statistico dell'Università di Firenze – Firenze.

It is well known that the simple linear correlation coefficient r admits two interesting geometrical interpretations as a cosine of an angle. The purpose of this paper is to give a geometric interpretation of the correlation ratio h along lines parallel to those now mentioned for r .

1.

Given the centered values \tilde{y}_i, \tilde{x}_i ($i=1,2,\dots,n$) assumed by the quantitative characteristics (variables) Y, X on n objects (individuals), it is well known that several geometric interpretations of the simple linear correlation coefficient r are possible. One of these consists in regarding \tilde{y}_i, \tilde{x}_i ($i=1,2,\dots,n$) as vectors, say $\tilde{\mathbf{y}}, \tilde{\mathbf{x}}$, in E^n and in looking at the cosine of the angle α they form. Alternatively, considering in E^2 the concentration ellipse

$$\mathbf{z}' \mathbf{R}^{-1} \mathbf{z} = 1, \quad \mathbf{R} = \begin{bmatrix} 1 & r \\ r & 1 \end{bmatrix}, \quad -1 < r < +1$$

it can be shown that

$$r = \begin{cases} \cos(\angle T_2 T_1 T_4) & \text{if } r \geq 0 \\ \cos(\angle T_1 T_2 T_3) & \text{if } r \leq 0 \end{cases}$$

* Based on a research partially funded by CNR.

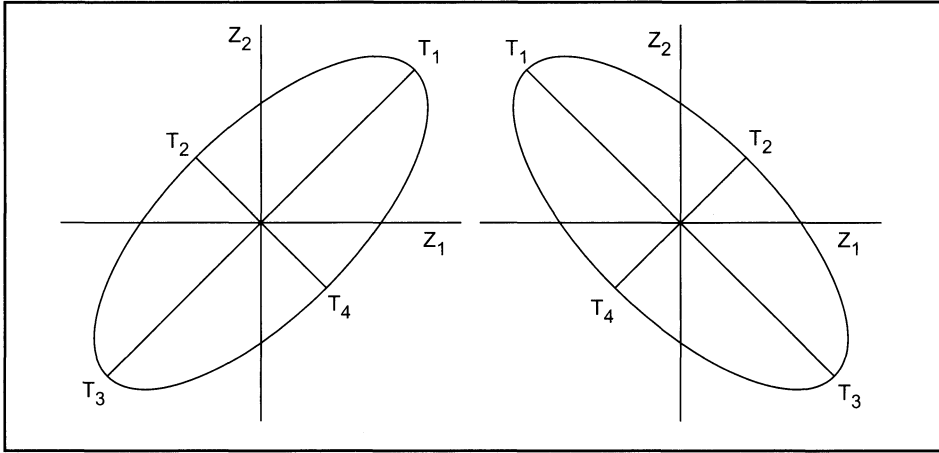


Fig. 1, 2.

where T_1, T_3 and T_2, T_4 are the intersections of the ellipse, respectively, with the major and minor axis (cf. Figures 1,2)¹.

But what about the correlation ratio η ? The purpose of this paper is to give a geometric interpretation of η along lines parallel to those mentioned above for r .

2.

Let us suppose that n objects are classified according to q values $y_1, \dots, y_k, \dots, y_q$ of a variable Y and to p attributes $a_1, \dots, a_h, \dots, a_p$ of a qualitative characteristic A . The results are usually summarized as in Table I, in which n_{hk} ($h=1, 2, \dots, p; k=1, 2, \dots, q$) denotes the number of objects taking the attribute a_h of A and the value y_k of Y and

$$n_{h\cdot} = \sum_{k=1}^q n_{hk} > 0, \quad n_{\cdot k} = \sum_{h=1}^p n_{hk} > 0, \quad \sum_{h=1}^p n_{h\cdot} = \sum_{k=1}^q n_{\cdot k} = n.$$

Setting

$$\bar{y}_h = \frac{\sum_{k=1}^q y_k n_{hk}}{n_{h\cdot}}, \quad \bar{y} = \frac{\sum_{k=1}^q y_k n_{\cdot k}}{n}$$

¹ In these interpretations it is assumed that E^n and E^2 are Euclidean spaces with scalar products represented – with respect to the natural bases of E^n and E^2 – by unit matrices of appropriate order.

Tab. I.

A			Y			Totals
	y_1	...	y_k	...	y_q	
a_i	n_{i1}	...	n_{ik}	...	n_{iq}	$n_{i\cdot}$
.
.
a_h	n_{h1}	...	n_{hk}	...	n_{hq}	$n_{h\cdot}$
.
.
a_p	n_{p1}	...	n_{pk}	...	n_{pq}	$n_{p\cdot}$
Totals	$n_{\cdot 1}$...	$n_{\cdot k}$...	$n_{\cdot q}$	n

the correlation ratio h is defined by the expression

$$\eta = \left(\frac{\sum_{h=1}^p (\bar{y}_h - \bar{y})^2 n_{h\cdot}}{\sum_{k=1}^q (y_k - \bar{y})^2 n_{\cdot k}} \right)^{1/2}$$

and interpreted as the square root of the ratio between the “interclass” variance and the total variance.

3.

In order to obtain a geometric interpretation of η , let us consider Table II where, besides Y , we have the new variables $X_1, \dots, X_h, \dots, X_p$. X_h ($h=1, 2, \dots, p$) denotes the variable “number of times with which the attribute a_h of A is present in an object”, and hence

$$x_{ih} = \text{value taken by } X_h \text{ on the } i\text{th } (i=1, 2, \dots, n) \text{ object}$$

$$= \begin{cases} 1 & \text{if the } i\text{th object takes the attribute } a_h \text{ of } A \\ 0 & \text{otherwise.} \end{cases}$$

In turn,

$$y_i = \text{value taken by } Y \text{ on the } i\text{th object.}$$

Tab. II.

Objects	X_1	...	X_h	...	X_p	Y
1	x_{11}	...	x_{1h}	...	x_{1p}	y_1
.
.
i	x_{i1}	...	x_{ih}	...	x_{ip}	y_i
.
.
n	x_{n1}	...	x_{nh}	...	x_{np}	y_n
Totals	$n_{1\cdot}$...	$n_{h\cdot}$...	$n_{p\cdot}$	$n\bar{y}$

Of course, since there is a one-to-one correspondence between Table I and Table II, their informative contents are the same.

Moreover, writing

$$\mathbf{X} = \begin{bmatrix} x_{11} & \dots & x_{1h} & \dots & x_{1p} \\ \cdot & \dots & \cdot & \dots & \cdot \\ x_{i1} & \dots & x_{ih} & \dots & x_{ip} \\ \cdot & \dots & \cdot & \dots & \cdot \\ x_{n1} & \dots & x_{nh} & \dots & x_{np} \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \cdot \\ y_i \\ \cdot \\ y_n \end{bmatrix}$$

we have

$$\mathbf{u}'\mathbf{X} = [n_{1\cdot} \dots n_{h\cdot} \dots n_{p\cdot}], \quad \mathbf{u}'\mathbf{y} = n\bar{y}$$

where $\mathbf{u}' = [1 \dots 1 \dots 1]$.

Therefore, the centered matrix and vector, corresponding to \mathbf{X} and \mathbf{y} , are

$$\tilde{\mathbf{X}} = \mathbf{X} - \mathbf{u} \frac{1}{n} \mathbf{u}'\mathbf{X} = \mathbf{X} - \mathbf{u} \left[\frac{n_{1\cdot}}{n} \dots \frac{n_{h\cdot}}{n} \dots \frac{n_{p\cdot}}{n} \right], \quad \tilde{\mathbf{y}} = \mathbf{y} - \mathbf{u}\bar{y}.$$

Notice that – since we have assumed $n_{h\cdot} > 0$ ($h=1,2,\dots,p$) – $r(\mathbf{X})=p \leq n$ and $r(\tilde{\mathbf{X}})=p - 1$; thus, $\tilde{\mathbf{X}}$ is not of full column rank.

4.

With these premises, we are in a position to give a geometric interpretation of η . Firstly, let

- $S(\mathbf{u})$: the vector space spanned by the vector \mathbf{u} ;
- $S(\mathbf{X})$: the vector space spanned by the column vectors of \mathbf{X} ;
- $S(\tilde{\mathbf{X}})$: the vector space spanned by the column vectors of $\tilde{\mathbf{X}}$;

and

$P(\mathbf{u})$: the orthogonal projection matrix onto $S(\mathbf{u})$;

$P(\mathbf{X})$: the orthogonal projection matrix onto $S(\mathbf{X})$;

$P(\tilde{\mathbf{X}})$: the orthogonal projection matrix onto $S(\tilde{\mathbf{X}})$.

Then, since $S(\mathbf{u})$ and $S(\tilde{\mathbf{X}})$ are orthogonal complements in $S(\mathbf{X})^{(2)}$, $\forall \mathbf{x} \in E^n$ we have

$$P_{(X)} \mathbf{x} = P_{(u)} \mathbf{x} + P_{(\tilde{X})} \mathbf{x}$$

and hence

$$P_{(\tilde{X})} = P_{(X)} - P_{(u)} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' - \mathbf{u}(\mathbf{u}'\mathbf{u})^{-1}\mathbf{u}'.$$

Therefore, the orthogonal projection $\hat{\mathbf{y}}$ of $\tilde{\mathbf{y}}$ onto $S(\tilde{\mathbf{X}})$ is given by

$$\hat{\mathbf{y}} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\tilde{\mathbf{y}} - \mathbf{u}(\mathbf{u}'\mathbf{u})^{-1}\mathbf{u}'\tilde{\mathbf{y}} = \mathbf{X} \begin{bmatrix} \bar{y}_1 - \bar{y} \\ \vdots \\ \bar{y}_h - \bar{y} \\ \vdots \\ \bar{y}_p - \bar{y} \end{bmatrix}.$$

Now, denoting by α the angle between the vectors $\tilde{\mathbf{y}}$ and $\hat{\mathbf{y}}$, it is very easy to show that $\cos^2\alpha$ is just η^2 .

In fact, since $(\mathbf{X}'\mathbf{X} = \text{diag}[n_1, \dots, n_h, \dots, n_p])$

$$\hat{\mathbf{y}}'\hat{\mathbf{y}} = [\bar{y}_1 - \bar{y} \dots \bar{y}_h - \bar{y} \dots \bar{y}_p - \bar{y}] \mathbf{X}'\mathbf{X} \begin{bmatrix} \bar{y}_1 - \bar{y} \\ \vdots \\ \bar{y}_h - \bar{y} \\ \vdots \\ \bar{y}_p - \bar{y} \end{bmatrix} = \sum_{h=1}^p (\bar{y}_h - \bar{y})^2 n_h.$$

² Actually, for $\forall a \in R$ and $\forall \mathbf{b} \in R^p$, we have

(i) $\mathbf{X} = \tilde{\mathbf{X}} + u \frac{1}{n} \mathbf{u}' \mathbf{X} \Rightarrow \mathbf{X} \mathbf{b} = \tilde{\mathbf{X}} \mathbf{b} + u \left(\frac{1}{n} [\mathbf{u}' \mathbf{X} \mathbf{b}] \right)$

(ii) $\mathbf{a} \mathbf{u}' \tilde{\mathbf{X}} \mathbf{b} = 0$;

then, $S(\mathbf{X}) = S(\mathbf{u}) \oplus S(\tilde{\mathbf{X}})$ and $S(\mathbf{u}) \perp S(\tilde{\mathbf{X}})$.

and

$$\tilde{\mathbf{y}}'\tilde{\mathbf{y}} = (\mathbf{y} - u\bar{\mathbf{y}})'(\mathbf{y} - u\bar{\mathbf{y}}) = \mathbf{y}'\mathbf{y} - \mathbf{u}'\mathbf{n}\bar{\mathbf{y}}^2 = \sum_{k=1}^q (y_k - \bar{y})^2 n_{\cdot k}$$

we have⁽³⁾

$$\cos^2 \alpha = \frac{\hat{\mathbf{y}}'\hat{\mathbf{y}}}{\tilde{\mathbf{y}}'\tilde{\mathbf{y}}} = \eta^2.$$

5.

Another geometric interpretation of η is easily obtained by considering in E^2 the ellipse

$$\mathbf{z}'\mathbf{N}^{-1}\mathbf{z} = 1, \quad \mathbf{N} = \begin{bmatrix} 1 & \eta \\ \eta & 1 \end{bmatrix}, \quad 0 \leq \eta < +1.$$

Denoting by d and D the lengths of the minor and major axes of the ellipse (cf. Figure 1), since

$$d = 2\sqrt{1-\eta}, \quad D = 2\sqrt{1+\eta}$$

we have

$$\eta = \frac{D^2 - d^2}{D^2 + d^2} = \frac{1 - \frac{d^2}{D^2}}{1 + \frac{d^2}{D^2}} = \frac{1 - t^2}{1 + t^2}.$$

But $t = d/D = \tan(T_2 T_1 T_4/2)$; thus,

$$\eta = \cos(T_2 T_1 T_4).$$

Of course, if the angle $T_2 T_1 T_4$ is rectangular then the correlation ratio is 0. In turn, narrow ellipses correspond to angles close to 0 and, therefore, to correlation ratios close to +1.

³ Notice that this is equivalent to the squared multiple linear correlation coefficient between Y and the set of variables $X_1, \dots, X_p, \dots, X_p$.

REFERENCES

- Châtillon G.: *The Balloon Rules for a Rough Estimate of the Correlation Coefficient*. The American Statistician, pp. 58–60; 1984.
- Châtillon G.: *Reply*. The American Statistician, pp. 330–331; 1984.
- Gypen L.M.J.: *Comment on Rogers and Nicewander*. The American Statistician, pp. 291, 1988.
- Marks E.: *A note on the Geometric Interpretation of the Correlation Coefficient*. Journal of Educational Statistics, pp. 233–237; 1982.
- Ozer D.J.: *Correlation and the Coefficient of Determination*. Psychological Bulletin, pp. 307–315; 1985.
- Rodgers J.L., Nicewander W.A.: *Thirteen Ways to Look at the Correlation Coefficient*. The American Statistician, pp. 59–66; 1988.
- Rodgers J.L.: *Reply*. The American Statistician, p. 291, 1988.
- Schilling M.F.: *Some Remarks on Quick Estimation of the Correlation Coefficient*. The American Statistician, p. 330, 1984.
- Thöni H.: *Comment on Rogers and Nicewander*. The American Statistician, p. 290, 1988.

SUMMARY

È ben noto che il coefficiente di correlazione lineare semplice r ammette due interessanti interpretazioni geometriche come coseno di un angolo. Lo scopo di questa nota è quello di mostrare la possibilità di interpretare in termini analoghi il rapporto di correlazione η .