

The Effect on R^2 of Adding an Independent Variate in
Multivariate Regression

Let the $(k+1)$ th independent variate be x_t , the k other variates being x_{ti} , $t=1,2,\dots, T$, $i=1,2,\dots,k$. Without loss of generality the x_{ti} are assumed to be orthogonal. Also w.l.g. the dependent variate y_t and all the independents are assumed to be standardized, i.e. their means are zero and their sum squares equal T . Let R_k^2 and R_{k+1}^2 be the values for the k and $(k+1)$ regressions respectively. The full $(k+1)$ regression is then

$$(1) \quad y_{tc} = cx_t + \sum_{i=1}^k x_{ti} b_i$$

Let

$$(2) \quad \frac{1}{T} \sum_t y_t x_t = p; \quad \frac{1}{T} \sum_t y_t x_{ti} = q_i; \quad \frac{1}{T} \sum_t x_t x_{ti} = r_i, \quad i=1,2,\dots,k$$

Obviously p and the q_i and r_i are correlation coefficients. They are, in fact, all the simple c.c.'s in the system since, by hypothesis, $\sum_t x_{ti} x_{tj}$ ($i, j=1,2,\dots,k$, $i \neq j$) is zero. Then

$$(3) \quad R_k^2 = \sum_{i=1}^k q_i^2$$

From (1) and (2),

$$(4) \quad R_{k+1}^2 = c^2 + 2c \sum b_i r_i + \sum b_i^2,$$

the normal equations being

$$(5) \quad p = c + \sum r_i b_i \\ q_i = r_i c + b_i, \quad i = 1,2,\dots,k.$$

On substitution for i and the b_i from (5) into (4), and on reduction,

$$(6) \quad R_{k+1}^2 = [(p - \sum q_i r_i)^2 + (1 - \sum r_i^2) \sum q_i^2] / (1 - \sum r_i^2),$$

or, using (3),

$$(7) \quad R_{k+1}^2 = R_k^2 + (p - \sum q_i r_i)^2 / (1 - \sum r_i^2)$$

The second term on the right of (7) has all the look of the square of a partial c.c. This turns out to be nearly the case. In fact, let

$$(8) \quad z_t = \sum_i x_{ti} r_i.$$

After some simple reduction we find

$$(9) \quad R_{k+1}^2 = R_k^2 + (1-r_{yz}^2) r^2(x,y/z)$$

with (see (2))

$$(10) \quad r_{xy} = p; \quad r_{xz} = \sqrt{\sum r_i^2}; \quad r_{yz} = \sum q_i r_i / \sqrt{\sum r_i^2}.$$

of course, $R_{k+1}^2 \geq R_k^2$.

No novelty is claimed for (7) or (9). The object of this note is to show a particular aspect of the fundamental role of orthogonization of independent variates in multivariate analysis, to say nothing of the simplification in exposition of standardization of all variates. Of course, a result like (7) can be produced using the raw data, but the apparent simplicity in that case of the second term on the right is spurious as consisting of a complicated matrix expression.

5 April 1967

R. C. Geary