# University of the Philippines SCHOOL OF ECONOMICS

Discussion Paper 8101

February 1981

RELATIVE CONTRIBUTIONS OF MIXED VARIABLES TO THE VARIATION OF A REGRESSAND

by

José Encarnación, Jr.

NOTE: UPSE Discussion Papers are preliminary versions circulated privately to elicit critical comment. They are 'protected by the Copyright Law (PD No. 49) and are not for quotation or reprinting without prior approval.

### Abstract

This note shows the direct comparability of the beta coefficients of ordinary scalar variables and of classificatory vector variables.

Accordingly, even when both kinds of variables appear in a regression equation, their relative contributions to the variation of the regressand can be ranked by their beta squares.

# Relative Contributions of Mixed Variables to the Variation of a Regressand

#### J. Encarnación

Consider a regression equation whose regressions include classificatory as well as ordinary scalar variables. A classificatory variable is essentially a vector that has as many components as there are different (mutually exclusive and exhaustive) categories in the classification. For example, one might estimate a regression equation that explains employees' salaries in terms of length of service (a scalar), occupation (a classificatory variable), etc. One might then want to estimate the relative contributions of the explanatory variables to the variation of the dependent variable. Handling this problem by beta coefficients is well known when the explanatory variables are all of one kind, either all scalar or all classificatory. There seems, however, to be no convenient reference that discusses this matter when the explanatory variables are mixed, i.e. when they include both kinds. This expository note might therefore be of some use.

1

Let  $\mathbf{x}=(\mathbf{x}_0,\,\mathbf{x}_1,\,\ldots,\,\mathbf{x}_K)$  where  $\mathbf{x}_k=1$  for an individual (or observation) if it belongs to category  $\mathbf{k}$  ( $\mathbf{k}=0,\,1,\,\ldots,\,K$ ) of classification  $\mathbf{x},\,\,\mathbf{x}_k=0$  otherwise, and  $\sum_{k=0}^K \mathbf{x}_k=1$ . More precisely, for any given individual  $\mathbf{i},\,\,\mathbf{x}_{ki}=1$  if  $\mathbf{i}$  is in category  $\mathbf{k},\,\,0$  otherwise, and  $\sum_{k=0}^K \mathbf{x}_{ki}=1$ . To each  $\mathbf{i}$  thus corresponds  $\mathbf{x}_i=(\mathbf{x}_{0i},\,\,\mathbf{x}_{1i},\,\,\ldots,\,\,\mathbf{x}_{Ki})$ .

Suppose it is appropriate to explain y in terms of x, z, u and v by means of a regression equation, where z is another classificatory variable  $(z_0, z_1, \ldots, z_J)$  while u and v are real variables. (Discussion of more than two variables of either kind would be straightforward.) We calculate

(1) 
$$y' = c + \sum_{j=1}^{K} a_{k}^{\alpha} x_{k} + \sum_{j=1}^{J} b_{j}^{\alpha} z_{j} + p(u - \bar{u}) + q(v - \bar{v})$$

where the  $a_k^k$ ,  $b_j^k$ , p and q are the regression coefficients and y' is the predicted y. As usual, overbars denote means. Note that  $x_0$  and  $z_0$  are omitted in (1) in order to have determinate coefficients (Suits 1957).

We want to express (1) in the form

(2) 
$$y' = \bar{y} + \sum_{0}^{K} a_k x_k + \sum_{0}^{J} b_j z_j + p(u - \bar{u}) + q(v - \bar{v})$$

where  $\mathbf{x}_0$  and  $\mathbf{z}_0$  are included, and the  $\mathbf{a}_k$  and  $\mathbf{b}_j$  measure the effects on an individual's y resulting from its belonging to k of x and to j of z, respectively. It is to be noted that the  $\mathbf{a}_k$  and  $\mathbf{b}_j$ , which might be called category effects (Encarnación 1975), are measured from  $\tilde{\mathbf{y}}$ . For suppose that for an individual i,  $\mathbf{x}_{ki}$  = 1 for a particular k and  $\mathbf{z}_{ji}$  = 1 for a particular j. Then

$$y'_{i} = \bar{y} + a_{k} + b_{j} + p(u_{i} - \bar{u}) + q(v_{i} - \bar{v})$$

so that  $a_k$  and  $b_i$  are simply added on to  $\bar{y}$ .

From least squares properties, using (1),

$$(3) \qquad \mathbf{c} = \bar{\mathbf{y}} - \sum_{1}^{K} \mathbf{a}_{\mathbf{k}}^{*} \, \bar{\mathbf{x}}_{\mathbf{k}} - \sum_{1}^{J} \mathbf{b}_{\mathbf{j}}^{*} \, \bar{\mathbf{z}}_{\mathbf{j}} - \mathbf{p}(\bar{\mathbf{u}} - \bar{\mathbf{u}}) - \mathbf{q}(\bar{\mathbf{v}} - \bar{\mathbf{v}})$$

$$= \bar{\mathbf{y}} - \sum_{1}^{K} \mathbf{a}_{\mathbf{k}}^{*} \, \bar{\mathbf{x}}_{\mathbf{k}} - \sum_{1}^{J} \mathbf{b}_{\mathbf{j}}^{*} \, \bar{\mathbf{z}}_{\mathbf{j}}.$$

But c is also the predicted y for an individual satisfying  $x_0 = 1$ ,  $z_0 = 1$ ,  $u = \bar{u}$  and  $v = \bar{v}$ . Therefore

(4) 
$$a_0 = -\sum_{1}^{K} a_k^* \bar{x}_k$$

(5) 
$$b_0 = -\int_1^J b_j^* \bar{z}_j$$
.

Further, if an individual satisfies  $x_k = 1$   $(k \neq 0)$ ,  $z_0 = 1$ ,  $u = \overline{u}$ ,  $v = \overline{v}$ , the predicted y is  $c + a_k^a$ . Since we already know from (3)-(5) that

(6) 
$$c = \bar{y} + a_0 + b_0$$

we have  $c + a_{k}^{*} = \bar{y} + (a_{0} + a_{k}^{*}) + b_{0}$  so that

(7) 
$$a_k = a_0 + a_k^*$$
  $k = 1, ..., K$ 

The b are similarly determined.

Substituting (6) in (1),

(8) 
$$y' = \bar{y} + a_0 + b_0 + \sum_{j=1}^{K} a_k^* x_k + \sum_{j=1}^{J} b_j^* z_j + p(u - \bar{u}) + q(v - \bar{v})$$
  

$$= \bar{y} + a_0 + b_0 + \sum_{j=1}^{K} (a_k - a_0) x_k + \sum_{j=1}^{J} (b_j - b_0) z_j + p(u - \bar{u}) + q(v - \bar{v})$$

$$+ q(v - \bar{v})$$

$$= \bar{y} + a_0 \left(1 - \sum_{1}^{K} x_k\right) + \sum_{1}^{K} a_k x_k + b_0 \left(1 - \sum_{1}^{J} z_j\right) + \sum_{1}^{J} b_j z_j$$

$$+ p(u - \bar{u}) + q(v - \bar{v}).$$

But 
$$1 - \sum_{1}^{K} x_{k} = x_{0}$$
 and  $1 - \sum_{1}^{J} x_{j} = x_{0}$ ; hence (2).

We note for later reference that  $\bar{x}_k = n_k / n$ , where  $n_k$  is the number of individuals for which  $x_{ki} = 1$  and n is the total number of individuals. Also, as one might expect,

i.e. the mean  $\sum_{0}^{K} a_k x_k = 0$  (in the same way that the mean  $p(u - \bar{u})$ , say, is zero). For, multiplying (7) by  $n_k$ , summing both sides and then adding  $n_0$ ,  $a_0$  to the results,

$$\sum_{0}^{K} n_{k} a_{k} = n a_{0} + \sum_{1}^{K} n_{k} a_{k}^{st}$$

which, in view of (4), gives (9).

II

The motivation for calculating the partial beta coefficients of standard multiple regression is to be able to compare the relative contributions of the explanatory (scalar) variables to the variation of the dependent variable (see, e.g., Ezekiel and Fox 1959, p.,196).

Accordingly, the variables are standardized to zero means and unit variances, so that their beta coefficients become directly comparable.

Similarly, the beta coefficients discussed by Morgan et al. (1962) perform the same function in the case of classificatory variables. Our problem

is to see whether all the beta coefficients in a regression with mixed variables are directly comparable.

Write

$$\frac{(10)}{s_y} = \beta_x f(x) + \beta_z q(z) + \beta_u \frac{u - \bar{u}}{s_u} + \beta_v \frac{v - \bar{v}}{s_v}$$

which is to be equivalent to (cf. (2))

(11) 
$$\frac{y' - \bar{y}}{s_y} = \frac{\sum_{0}^{K} a_k x_k}{s_y} + \frac{\sum_{0}^{J} b_j z_j}{s_y} + \frac{p(u - \bar{u})}{s_y} + \frac{q(v - \bar{v})}{s_y}$$

where s is the standard direction of y, etc.,

(12) 
$$\beta_u = p s_u/s_y$$

which is the textbook definition of a partial beta coefficient, similarly for  $\beta_{_{\boldsymbol{V}}},$ 

(13) 
$$\beta_{x} = \frac{(\sum_{0}^{K} a_{k}^{2} n_{k} /(n - 1))^{1/2}}{s_{v}}$$

from Morgan et al. (1962), and the functions f(x) and g(z) are implicitly defined by the equivalence of (10) and (11) and the definitions of the  $\beta$ 's. It is clear that if  $\frac{1}{u} > \beta_v^2$ , u contributes more than does v to the explanation of y variation. Our object is to show that f(x), say, standardizes x essentially in the same way that  $(u - \tilde{u})/s_u$  standardizes u, so that all the beta coefficients are then directly comparable.

From (10), (11) and (13), for individual i,

(14) 
$$f(x_i) = \frac{\sum_{k=0}^{K} a_k x_{ki}}{(\sum_{k=0}^{K} a_k^2 n_k /(n-1))^{1/2}}$$

from which

(15) 
$$f(x_i)^2 = \frac{\sum_{k=0}^{K} a_k^2 x_{ki}^2}{\sum_{h=1}^{n} \sum_{k=0}^{K} a_k^2 x_{kh}^2 / (n-1)}$$

since cross-product terms vanish and  $x_{ki} = x_{ki}^2$  (because  $x_{ki} = 0$  or 1 and  $\sum_{k=0}^{K} x_{ki} = 1$ ). But

$$\frac{(u_i - \bar{u})^2}{s_u^2} = \frac{p^2(u_i - \bar{u})^2}{\sum_{h=1}^n p^2(u_h - \bar{u})^2/(n-1)}$$

corresponds precisely to (15), the only difference being that while one can factor out  $p^2$  in (16), which of course does not affect the ratio, it is not possible to factor out  $\sum_{0}^{k} a_k^2$  in (15), which pertains to a vector. The key observation is that x being a classificatory variable,  $\sum_{k=0}^{K} a_k x_{ki}$  is the analogue of  $p(u_i - \bar{u})$  and both have zero means.

This completes our task, and all the beta squares may then be ranked to indicate the relative contributions of their corresponding variables to the explanation of y variation.

#### References

- Encarnación, J. "Income Distribution in the Philippines: The Employed and the Self-Employed," in <u>Income Distribution</u>, <u>Employment and Economic Development in Southeast and East Asia</u> (Tokyo: JERC and Manila: CAMS), 1975, pp. 742-775.
- Ezekiel, M. and Fox, K.A. Methods of Correlation and Regression Analysis, 3rd ed. (New York: Wiley), 1959.
- Morgan, J.N. et al. Income and Welfare in the United States (New York: McGraw-Hill), 1962. Appendix E.
- Suits, D.B. "Use of Dummy Variables in Regression Equations," <u>Journal</u> of the American Statistical Association, December 1957, pp. 548-551.