

Institute of Economic Development and Research

SCHOOL OF ECONOMICS
University of the Philippines System

Discussion Paper No. 74-3

April 19, 1974

ON MEASURING THE RELATIVE CONTRIBUTIONS OF
SEVERAL CLASSIFICATORY VARIABLES IN THE EXPLANATION
OF A DEPENDENT VARIABLE

by

José Encarnación, Jr.

Note: IEDR Discussion Papers are preliminary versions circulated privately to elicit critical comment. References in publications to Discussion Papers should be cleared with the author.

ON MEASURING THE RELATIVE CONTRIBUTIONS OF
SEVERAL CLASSIFICATORY VARIABLES IN THE EXPLANATION
OF A DEPENDENT VARIABLE

by J. Encarnación

Consider the explanation of a variable y (the logarithm of income for example) in terms of two classificatory variables A and B . A could be, say, an industrial sector classification with different categories such as manufacturing, agriculture, etc. Each classification consists of categories which are exhaustive and mutually exclusive, and we wish to explain an individual's y as a result of his belonging to some category of A and some category of B . We also want to estimate the relative contributions of those variables to the explanation of y variation. The purpose of this note is to make more explicit the treatment of these matters by Suits (1957), Morgan et al. (1962) and Oey (1972). It will be apparent how the discussion would proceed if there were three or more classificatory variables.

Let the categories of A be indexed by k ($k = 0, 1, \dots, K$), those of B by j ($j = 0, 1, \dots, J$); and let y_{kji} be the i th observation in the (k, j) -cell, $i = 1, \dots, n_{kj}$. There are n observations, so

$$n = \sum_{k=0}^K \sum_{j=0}^J n_{kj}$$

Write

$$n_{k.} = \sum_{j=0}^J n_{kj}$$

for the number of observations in the kth category of A; similarly,

$$n_{.j} = \sum_{k=0}^K n_{kj}$$

The mean y in category k of A is given by

$$\bar{y}_k = \frac{\sum_{j=0}^J \sum_{i=1}^{n_{kj}} y_{kji}}{n_k}$$

and for the mean y we write \bar{y} .

Let

$$x_k = \begin{cases} 1 & \text{if an observation belongs to } k \text{ of } A \\ 0 & \text{otherwise} \end{cases}$$

and define z_j similarly. With these dummy variables we can calculate a regression equation

$$(1) \quad y' = c + a'_1 x_1 + \dots + a'_K x_K + b'_1 z_1 + \dots + b'_J z_J$$

Note that the variables x_0 and z_0 are omitted from the regression in order to get determinate estimates of the coefficients; see Suits (1957). (Cf. the fact that a sex classification has two categories - male and female - but only one dummy variable would be included in a regression equation, i.e. one dummy variable is omitted.) For a given individual, at most one of the x_k , $k = 1, \dots, K$, and at most one of the z_j , $j = 1, \dots, J$, can be nonzero. Eq. (1) thus shows y' , the predicted y , to depend on an individual's membership in the different categories of A and of B.

We wish to write (1) in the form

$$(2) \quad y' = \bar{y} + a_k + b_j$$

where $k = 0, 1, \dots, K$ and $j = 0, 1, \dots, J$. Then a_k would measure the effect on y resulting from belonging to category k of A , which effect is measured from \bar{y} . Similarly for b_j in regard to j of B .

In order to get the a_k , from least squares regression properties we know that in (1),

$$\begin{aligned} c &= \bar{y} - \sum_1^K a'_k \bar{x}_k - \sum_1^J b'_j \bar{z}_j \\ &= \bar{y} - \sum_1^K a'_k n_{.k} / n - \sum_1^J b'_j n_{.j} / n \end{aligned}$$

But c is the predicted y when one belongs to category 0 of A (in which case $x_1 = \dots = x_K = 0$) and 0 of B . From this it follows that

$$(3) \quad \begin{aligned} a_0 &= - \sum_1^K a'_k n_{.k} / n \\ b_0 &= - \sum_1^J b'_j n_{.j} / n \end{aligned}$$

On the other hand, if one belongs to category k ($k = 1, \dots, K$) of A and 0 of B , the predicted y is $c + a'_k$. Accordingly we have

$$(4) \quad a_k = a_0 + a'_k \quad (k = 1, \dots, K)$$

Eqs. (3) and (4) thus determine the a_k , and a similar procedure gives the b_j .

Multiplying (4) by n_k , summing both sides and then adding $n_0 a_0$ to the results, one gets

$$\sum_0^K n_k a_k = n a_0 + \sum_1^K n_k a'_k$$

which, in view of (3), implies that the weighted sum of the a_k is zero:

$$(5) \quad \sum_0^K a_k n_k / n = 0$$

In order now to get (2) from (1), write

$$\begin{aligned} y' &= \bar{y} + a_0 + b_0 + \sum_1^K a'_k x_k + \sum_1^J b'_j z_j \\ &= \bar{y} + a_0 + b_0 + \sum_1^K (a_k - a_0) x_k + \sum_1^J (b_j - b_0) z_j \\ &= \bar{y} + a_0 \left(1 - \sum_1^K x_k\right) + \sum_1^K a_k x_k + b_0 \left(1 - \sum_1^J z_j\right) + \sum_1^J b_j z_j \end{aligned}$$

But $1 - \sum_1^K x_k = x_0$ and similarly for the other set of dummy variables, and we obtain (2). Hence (1) and (2) yield identical predicted values (and therefore identical error terms).

Define a new variable

$$y_{kji}^* = y_{kji} - a_k - b_j$$

by removing from each observation the effects of the classificatory variables as calculated. Then

$$y^* - \bar{y} = y - y^*$$

becomes a measure of the unexplained error.

The beta coefficients discussed by Morgan et al. serve to measure the relative contributions of the classificatory variables to the explanation of y variation:

$$(6) \quad \beta_A = \frac{(\sum_{k=0}^K n_k \cdot a_k^2 / (n - 1))^{1/2}}{s_y}$$

where s_y is the standard deviation of y . Similarly

$$\beta_B = \frac{(\sum_{j=0}^J n_j \cdot b_j^2 / (n - 1))^{1/2}}{s_y}$$

If, say, $\beta_A > \beta_B$, that would indicate that more of the variation in y is due to the classificatory variable A . As Morgan et al. point out, these coefficients are analogous to the partial beta coefficients of standard multiple regression.

Following Oey (1972), an approximate F-test could be used to test the significance of a classificatory variable in explaining y . For from (2),

$$(7) \quad e_{kji} = y_{kji} - \bar{y} - a_k - b_j = y_{kji}^* - \bar{y}$$

where e is the error term, and under certain assumptions, $(n - 1)s_e^2/\sigma_y^2$ would be a chi-square variable with $n - K - J - 1$ degrees of freedom. Accordingly in the case of A ,

$$(8) \quad F_A = \frac{\sum_{k=0}^K n_k \cdot a_k^2 / K}{\sum_{k=0}^K \sum_{j=0}^J \sum_{i=1}^{n_{kj}} (y_{kji}^* - \bar{y})^2 / (n - K - J - 1)}$$

If the sum of squares in the numerator is small relative to that of the denominator, that would indicate that the use of A does not add much to the explanation of y .

Finally, we consider the interesting (but false) conjecture that

$$R^2 = \beta_A^2 + \beta_B^2$$

where $R^2 = 1 - s_e^2/s_y^2$ is the coefficient of determination given by

(1). One might argue in the following way. We have

$$(9) \quad y_{kji} - \bar{y} = a_k + b_j + (\bar{y}_{kj} - \bar{y} - a_k - b_j) + (y_{kji} - \bar{y}_{kj})$$

where \bar{y}_{kj} is the mean y in the (k, j) -cell, so that referring to (7),

$$(10) \quad e_{kji} = (\bar{y}_{kj} - \bar{y} - a_k - b_j) + (y_{kji} - \bar{y}_{kj})$$

i.e. the error term is the sum of two differences: between the cell mean and the predicted value given by (1) or (2), and between the actual y and the cell mean. From (9),

$$(11) \quad \sum \sum \sum_{kji} (y - \bar{y})^2 = \sum \sum \sum [a_k + b_j + (\bar{y}_{kj} - \bar{y} - a_k - b_j) + (y_{kji} - \bar{y}_{kj})]^2$$

Suppose cross-product terms on the right-hand side all vanish. Then

$$\sum \sum \sum_{kji} (y - \bar{y})^2 = \sum \sum \sum a_k^2 + \sum \sum \sum b_j^2 + \sum \sum \sum (\bar{y}_{kj} - \bar{y} - a_k - b_j)^2 + \sum \sum \sum (y - \bar{y}_{kj})^2$$

in which case

$$(*) \quad \frac{(n-1)s_y^2}{\sigma_y^2} = \frac{(n-1)s_a^2}{\sigma_y^2} + \frac{(n-1)s_b^2}{\sigma_y^2} + \frac{(n-1)s_p^2}{\sigma_y^2} + \frac{(n-1)s_q^2}{\sigma_y^2}$$

where $s_a^2 = \sum_k n_k \cdot a_k^2 / (n-1)$, $s_p^2 = \sum \sum \sum (\bar{y}_{kj} - \bar{y} - a_k - b_j)^2 / (n-1)$, etc. Assuming that y is normally distributed, it is known that $(n-1)s_y^2 / \sigma_y^2$ is a chi-square variable with $n-1$ degrees of freedom, so that the terms on the right-hand side of (*) must then be independent chi-square variables whose degrees of freedom add up to $n-1$. That is,

$$(**) \quad \chi_{y,n-1}^2 = \chi_{a,K}^2 + \chi_{b,J}^2 + \chi_{p,KJ}^2 + \chi_{q,n-KJ-K-J-1}^2$$

One could then get the F-statistic (8) by adding, in view of (10), the degrees of freedom in the last two terms of (**). In other words, e would also be chi-square with $n - K - J - 1$ degrees of freedom, and

$$(***) \quad s_e^2 = s_p^2 + s_q^2$$

Then from (*) and (***) one has

$$\begin{aligned} 1 &= \frac{s_a^2}{s_y^2} + \frac{s_b^2}{s_y^2} + \frac{s_e^2}{s_y^2} \\ &= \beta_A^2 + \beta_B^2 + \frac{s_e^2}{s_y^2} \end{aligned}$$

which says that β^2 is the sum of the β^2 s. This statement is wrong, however, as the cross-product terms on the right-hand side of (11) do not all vanish. For instance,

$$\begin{aligned} \sum \sum a_k (\bar{y}_{kj} - \bar{y} - a_k - b_j) &= \sum_k a_k n_k \bar{y}_k - \bar{y} \sum_k n_k a_k \\ &\quad - \sum_k n_k a_k^2 - \sum \sum a_k b_j \end{aligned}$$

The second and fourth terms on the right-hand side are easily seen to vanish, but

$$\sum_k a_k n_k \bar{y}_k - \sum_k n_k a_k^2 = \sum_k a_k n_k (\bar{y}_k - a_k)$$

is zero only if $\bar{y}_k - a_k$ is a constant, as would be the case if $a_k = \bar{y}_k - \bar{y}$. The last clearly does not hold in general, however. If it did, the calculated a_k would not change as one adds more classificatory variables in the explanation of y , but in fact the a_k - which derive from the regression coefficients of (1) - would change as we add more variables in the specification.

The author is indebted to Professor Lydia H. Flores, of the University of the Philippines Department of Mathematics, for very helpful discussions on the subject of this note.

References

- Hogg, R.V. and Craig, A.T. Introduction to Mathematical Statistics, 2nd ed., Macmillan, 1965.
- Morgan, J.N. et al. Income and Welfare in the United States, McGraw-Hill, 1962. Appendix E.
- Oey, M. "Income Distribution in Brazil," unpublished Ph.D. dissertation, University of California at Berkeley, 1972.
- Scheffé, H. The Analysis of Variance, Wiley, 1959.
- Suits, D.R. "Use of Dummy Variables in Regression Equations," J. Amer. Stat. Assn., December 1957, pp. 548-51.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100