# DECOMPOSING THE COEFFICIENT OF DETERMINATION IN MULTIPLE REGRESSION

## LEE M. WOLFLE

In multiple regression, the coefficient of determination (or R-square) has a very useful interpretation. The statistic is the ratio of the variation that is explained by the regression equation to the total variation of the dependent variable. For example, a coefficient of determination equal to .45 indicates that the independent variables can explain 45 percent of the variation in the dependent variable.

It follows immediately that a person's next instinct is to want to allocate among the several independent variables the explained variation in the dependent variable. For example, many people would like to say that if the regression of a dependent variable on three independent variables explains 45 percent of the variance, that (say) 25 percent was due to the first independent variable, 15 percent to the second, and 5 percent to the third. While this interpretation is tempting, it should be avoided. The reason for avoiding it is that there is no unique way of decomposing the explained variance, and if there is no unique way of doing so, then there is no meaningful way of doing so.

For example, consider the regression of one dependent variable, Y, on two independent variables:

$$Y = a + bX + cZ + u,$$

1

where Y is the dependent variable, X and Z are independent variables, "a" is the intercept, "b" and "c" are partial regression coefficients, and "u" is the disturbance or error term. From this regression, one would obtain the coefficient of determination, $R^2y.xz$, which is merely a convenient notation for R-square of Y regressed on X and Z. By the method of part correlations, it may be shown that:

$$R^2y.xz = R^2y.x + R^2y(z.x)$$

where $R^2y.x$ is the squared zero-order correlation of Y and X, and $R^2y(z.x)$ is the squared part correlation of Y with Z residualized for X. (This equation is analogous to formula 5.10 in Kerlinger and Pedhazur, 1973.)

But it is also true that:

$$R^2y.xz = R^2y.z + R^2y(x.z),$$

and in general $R^2y.x$ does not equal $R^2y(x.z)$. If these two quantities are not equal, by which quantity therefore does one measure the unique contribution of X to the explained variation in Y? Because there are two answers, two different answers, there is no unique solution.

These quantities may appear mysterious in symbolic form, but they are familiar quantities, which appear in the SPSS regression output. They appear in the summary table in a column of numbers entitled "R-Square Change." People often want to interpret these quantities as measuring the amount of variance explained by each independent variable. This temptation should be avoided.

Suppose, for example, that one regresses educational attainment on two independent variables, father's education and father's occupational

status. Suppose the coefficient of determination for this equation is equal to .31. If father's education was added to the regression as the first independent variable, then one would learn from the SPSS summary table that the R-Square Change for father's education was .27 and for father's occupation was .04. If, however, father's occupation were to have been listed first, then one would learn that the R-Square Change for father's occupation was .23 and for father's education was .08. Question: Does father's education explain 27 percent of variation in son's education, or does it explain 8 percent? The answer is, "Yes, it does." That is, without a unique way to decompose the explained variance, there is no unique answer to the question.

Let us try another approach. It is well known (e.g., Kerlinger and Pedhazur, 1973, formula 4.17) that the coefficient of determination may be decomposed into the sum of the products of the zero-order correlations and their associated beta-weights. One may, therefore, be tempted to interpret the product of say $R_{yx}$ times $B_{yx.z}$ (where $R_{yx}$ is the zero-order correlation of Y and X, and $B_{yx.z}$ is the beta-weight of Y regressed on X controlling for Z) as the amount of variance in Y explained by X. The problem with this approach, however, is that the zero-order correlation and the beta-weight are not constrained to have the same sign. In such cases, one would have to interpret the product as being a negative component to the explained variance, which is a very troublesome concept.

For example, consider the regression of son's occupational status on his educational attainment, his father's education, and his father's

occupational status. The correlations of son's occupation with these three independent variables are in one sample, respectively, .47, .38, and .73. The corresponding beta-weights are, respectively, .21, -.11, and .69. The coefficient of determination may be decomposed:

$$R^2 = .56 = (.21)(.47) + (-.11)(.38) + (.69)(.73)\ \text{or}$$

$$R^2 = .56 = .10 - .04 + .50.$$

It is with no relish whatsoever that one should interpret the amount of variance explained by father's education as being minus 4 percent. That is equivalent to saying that the addition of father's education to the regression equation takes away four percent of the variance in son's occupational status. Not only does that interpretation not make any sense in a substantive way, it is mathematically impossible.

Finally, consider the decomposition:

$$R^2y.xz = B^2y.x + B^2y.z + 2(By.x)(By.z)(Rxz),$$

In which $B^2y.x$ is the square of the beta-weight of Y regressed on X controlling for Z (and analogously for $B^2y.z$), By.x is the beta-weight of Y regressed on X, and Rxz is the zero-order correlation between X and Z. This decomposition seemingly contains a portion (the squared beta-weight) that can be uniquely attributed to the independent variable, but the decomposition also contains an explicit term (or more than one term if there are more than two independent variables) representing the contribution to the explained variance in Y that is shared by both independent variables. What this decomposition indicates is that the explained variance in Y cannot be decomposed into unique separate components due to each independent variable (unless Rxz = 0, a very rare occurrence).

If the coefficient of determination cannot be uniquely partitioned into amounts of variance explained by each independent variable, how then does one measure the contribution of each independent variable to the dependent variable?  The solution would appear to be to use either:

(1)  the metric partial regression coefficients; or

(2)  the standardized partial regression coefficients,

or beta-weights.

Notice that neither of these are interpretable as components of explained variance.  Metric partial regression coefficients are to be interpreted as the amount Y changes for a one-unit increase in one independent variable while the other independent variables are held constant. Standardized partial regression coefficients are interpretable as the number of standard deviations Y changes for a one standard deviation increase in one independent variable while the other independent variables are held constant.

The standardized regression coefficients have the advantage of being standardized.  That is, the size of the metric regression coefficients depend upon the metric in which the independent variables have been measured.  If one of the independent variables is income, for example, the metric regression coefficients will be different if income is measured in increments of thousand dollars versus increments of single dollars.  In any event, these coefficients will be different from those of another independent variable measured in, say, years of schooling. Standardized coefficients get around this problem by measuring all the variables in standard deviation units.  Thus, the standardized

coefficients are comparable among independent variables. A beta-weight of .5 for one independent variable means that Y changes twice as much as it does when another variable, which has a beta-weight of .25, changes.

Seemingly, therefore, the standardized regression coefficient is better than the metric coefficient. But wait! The standardized coefficient has to be standardized in terms of something, and that "something" turns out to be a quantity which is not invariant across either samples or populations. I am referring to the ratio of the standard deviations of the independent to the dependent variables. That is,

BETA = B (Sx/Sy),

where BETA is a beta-weight and B is the corresponding metric regression coefficient. If, for example, one is interested in comparing the effects of one independent variable on a dependent variable, and wants to compare the size of this effect across two populations (e.g., freshman versus sophomores, men versus women, blacks versus whites, etc.) then the beta-weights can change as a function of a change in the ratio of standard deviations; even while the structural coefficients, the metric coefficients, remain constant across populations.

Therefore, in reporting regression results one should always report both the standardized and the metric coefficients. The former are useful in comparing the relative effects of independent variables within a sample or population, while the latter are useful for comparing the relative effects of independent variables across samples or populations.

In any event, the amount of variance explained by each independent variable is not a quantity that can be uniquely estimated; the use of such estimates is to be discouraged.

(Kerlinger and Pedhazur, 1973, pp. 297-305, discuss a method called "Commonality Analysis," which can sometimes be used to estimate amounts of variance explained by each independent variable. In essence, the approach measures the portion of explained variance for a single independent variable as that portion unexplained by all of the other independent variables. The method results in measures of unique contributions and common contributions. In most real-life cases, the common portions far outweigh the unique portions. Another problem is the proliferation of higher-order commonalities. With five independent variables, commonality analysis produces five unique components, and 26 common components. In my view, metric and standardized regression coefficients are to be preferred in reporting the results of regression analyses.)

## REFERENCES

Kerlinger, Fred N., and Pedazur, Elazar J. Multiple Regression in Behavioral Research. New York: Holt, Rinehart and Winston, 1973.