

DECOMPOSITION OF R^2 IN MULTIPLE REGRESSION WITH CORRELATED REGRESSORS

Abraham Genizi

Agricultural Research Organization

Abstract: Let correlated regressors and a dependent variable have a joint distribution, and assume that a suitable regression model has been found. A decomposition of R^2 into components corresponding to the regressors is proposed. The components are to serve as descriptive intuitive statistics indicating the relative importance of each regressor with respect to its overall effect on the dependent variable. When it is possible to partition the set of regressors into mutually orthogonal subsets, the sum of components of this decomposition in any subset is equal to the multiple R^2 due to that subset. Each component consists of a subcomponent "specific" to that regressor, and of "common" subcomponents with each of the other regressors. A measure of deviation of the set of regressors from orthogonality is given to help in assessing the amount of approximation used in the decomposition.

Key words and phrases: Multiple correlation coefficient, Decomposition of explained variance, Relative importance of regressors, Geometrical representation, Measure of deviation from orthogonality.

1. Introduction

The question of the "percentage of variance of the dependent variable 'explained' by each of the explanatory variables" is frequently raised and discussed in scientific papers of various disciplines. This is a relevant and important issue for many cases in which a model is investigated after its construction has been completed successfully. However, no satisfactory solution has been found to this problem for the most common case, in which the explanatory variables are not orthogonal to each other. Moreover, it is well known (see e.g. Williams (1978)), that in such cases no decomposition of R^2 exists with a meaningful allocation of the overall regression effect to the individual regressors. Kruskal (1984), and Kruskal and Majors (1989) discussed the philosophical problems involved in an attempt to decompose a joint effect of two or more interdependent factors into its components, and reviewed several measures given in the literature for the "relative importance" of variables in various contexts. In the present work we outline the type of problems for which it may be justified to decompose R^2 into

components associated with the regressors; then, a decomposition will be presented with some desirable properties. Due to these properties we advocate the components to serve as approximate measures of the relative importance of the regressors in a given regression.

The regression problems for which a decomposition may be relevant are those in which the uncontrolled regressors and the dependent variable have a joint distribution, because the 'contribution' of a regressor that is controlled depends on the range over which it is varied. This restriction may be somewhat relaxed to include regressors which may be controlled, in cases where it may be expected that the ranges of the regressors, and in fact the whole covariance matrix of the regressors, will remain unchanged in the future. Actually, this may be the justification for presenting R^2 itself in many cases. Also, if the effect of a regressor on the dependent variable is through a second regressor, then interest would usually be on the marginal effect of the second regressor after accounting for the first. In general, no cause-effect interrelationships between regressors are assumed, because our aim is a simultaneous rather than sequential decomposition.

We are looking for decompositions after a satisfactory model has been found and only within the framework of this model. The component of R^2 associated with a regressor may change as a result of any change in the set of regressors. Thus, no comparisons of components will be valid across different models, and components will be inappropriate for variable selection purposes. Also, the interpretation of the components will be narrower than in the orthogonal case, e.g. the component of variable x_i will not necessarily be the deduction from R^2 obtained by holding x_i fixed. The components will satisfy some reasonable criteria and will be accompanied by a measure of deviation of the regressor set from orthogonality. They should be considered as descriptive intuitive statistics indicating the relative importance of each of the regressors with respect to their overall effect on the dependent variable.

2. The Model and a Criterion for Decomposition

Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$ and \mathbf{y} be $p + 1$ linearly independent vectors of n observations on p explanatory variables x_1, x_2, \dots, x_p and the dependent variable y , which have a joint distribution. Since we are interested only in correlations, we assume, without loss of generality, that all $p + 1$ vectors are standardized with 0 means. The vector of least squares solution, \mathbf{b} , to the model $E(\mathbf{y}) = \mathbf{X}\beta$ is $\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$, where $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p)$. The vector fitted to the observed \mathbf{x}_i 's, $\hat{\mathbf{y}} = \mathbf{X}\mathbf{b} = \sum_i \mathbf{x}_i b_i$, is well known to be the projection of \mathbf{y} onto the subspace $V(\mathbf{X})$ of R^n spanned by $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$. The multiple correlation coefficient estimate R is the cosine of angle between \mathbf{y} and $\hat{\mathbf{y}}$, and its square $R^2 = R^2_{y(\mathbf{X})} = \hat{\mathbf{y}}'\hat{\mathbf{y}}/\mathbf{y}'\mathbf{y}$.

Naturally, the components of any decomposition in the case of uncorrelated (orthogonal) regressors should be proportional to the squared (simple) correlation coefficients between \mathbf{y} and the corresponding regressors. Moreover it will required to satisfy the following

Criterion for orthogonal compatibility (OC): for a set of regressors partitioned into mutually uncorrelated subsets, the sum of components in any subset should be equal to the multiple R^2 between \mathbf{y} and that subset of regressors.

The term ‘component of a subset’ will be used for the sum of components of the regressors in that subset of regressors. For a decomposition satisfying Criterion OC it also follows in the situation of mutually uncorrelated subsets, that a nonsingular linear transformation on any subset will change neither the component of that subset nor those of the regressors in any other subset.

Three simple statistics satisfy Criterion OC if all \mathbf{x}_i 's are orthogonal but not in general. They are

- (i) the squared (simple) correlation coefficient $r_{\mathbf{y}\mathbf{x}_i}^2 = (\mathbf{y}'\mathbf{x}_i)^2 / (\mathbf{y}'\mathbf{y} \cdot \mathbf{x}_i'\mathbf{x}_i)$;
- (ii) the marginal reduction of $\mathbf{y}'\mathbf{y}$ due to \mathbf{x}_i following reductions due to all other \mathbf{x}_j 's ($j \neq i$), $R_{\mathbf{y}\mathbf{x}_i \cdot \mathbf{X}_{-i}}^2 = R^2 - R_{\mathbf{y}(\mathbf{X}_{-i})}^2$, where \mathbf{X}_{-i} is the matrix \mathbf{X} with column \mathbf{x}_i deleted; and
- (iii) the squared length of the vector component ($\hat{\mathbf{y}}_i$) of \mathbf{y} in the direction of \mathbf{x}_i divided by $\mathbf{y}'\mathbf{y}$, that is, $\hat{\mathbf{y}}_i'\hat{\mathbf{y}}_i / \mathbf{y}'\mathbf{y} = b_i^2 \mathbf{x}_i'\mathbf{x}_i / \mathbf{y}'\mathbf{y}$ which is proportional to the squared standardized partial regression coefficient.

Indeed, if the subspace of the regressors \mathbf{x}_1 and \mathbf{x}_2 is orthogonal to the other regressors, but \mathbf{x}_1 is not orthogonal to \mathbf{x}_2 , it can be seen easily that the sum of the two components for any of these three statistics in the subspace $V(\mathbf{x}_1, \mathbf{x}_2)$, is either larger or smaller than R^2 between \mathbf{y} and this subspace, which is $\hat{\mathbf{y}}_{12}'\hat{\mathbf{y}}_{12} / \mathbf{y}'\mathbf{y}$, where $\hat{\mathbf{y}}_{12}$ is the projection of \mathbf{y} into this subspace.

Another simple statistic is presented sometimes as the decomposition of R^2 . It is proportional to

- (iv) $\hat{\mathbf{y}}'\hat{\mathbf{y}}_i = b_i \mathbf{x}_i'\mathbf{y}$, or equivalently to the (signed) length of the projection of $\hat{\mathbf{y}}_i$ on $\hat{\mathbf{y}}$. Although it satisfies Criterion OC, it is not suitable to serve as a decomposition of R^2 , because some of its components in specific situations may turn out to be negative. This can be illustrated in two dimensions, in some of the cases when the projection $\hat{\mathbf{y}}_{12}$ of \mathbf{y} on $V(\mathbf{x}_1, \mathbf{x}_2)$ falls outside the acute angles, between \mathbf{x}_1 and \mathbf{x}_2 . Thus, although both $r_{\mathbf{y}\mathbf{x}_1}$ and $r_{\mathbf{y}\mathbf{x}_2}$ may be positive, the regression coefficient of the smaller r , say of \mathbf{x}_2 and consequently $\hat{\mathbf{y}}_2'\hat{\mathbf{y}}_{12}$, will be negative. They are said to serve as ‘‘corrections’’ to the too high values given to the corresponding coefficients of \mathbf{x}_1 . Pratt (1987) proposed $\hat{\mathbf{y}}'\hat{\mathbf{y}}_i$ as a measure of relative importance of the regressors. In his opinion, the possibility of negative values is ‘‘not a defect in definition’’, but it ‘‘signifies a situation too complex

for a single measure". He derived this measure by an axiomatic approach and gave to it several interpretations. (See the Discussion for comparisons with the measure proposed here.)

There is another way to define a decomposition of R^2 , which satisfies Criterion OC (see also Kruskal (1984), Pratt (1987)). Its component for x_i is

$$(v) \quad \bar{r}_i^2 = \frac{1}{p} \sum_{k=0}^{p-1} \binom{p-1}{k}^{-1} \sum_{\mathbf{X}_k} R_{yx_i.\mathbf{X}_k}^2, \quad i = 1, 2, \dots, p, \quad (1)$$

where the second summation goes over all possible groups of k variables, \mathbf{X}_k , out of x_1, x_2, \dots, x_p excluding x_i , $R_{yx_i.\mathbf{X}_k}^2 = R_{y(x_i.\mathbf{X}_k)}^2 - R_{y(\mathbf{X}_k)}^2$ and $R_{yx_i.\mathbf{X}_0}^2 = r_{yx_i}^2$. This is derived easily if we take all possible 'hierarchical' or stepwise decompositions of R^2 into elements of type $R_{yx_i.\mathbf{X}_k}^2$. (For example, for $p = 3$, six decompositions of the type $R^2 = R_{yx_1.x_2x_3}^2 + R_{yx_2.x_3}^2 + R_{yx_3}^2$ are obtained by permutation of the indices 1, 2 and 3.) Each decomposition contains p components, with exactly one component of the type $R_{yx_i.\mathbf{X}_k}^2$ for each i ($1 \leq i \leq p$). Averaging these hierarchical decompositions separately for elements of type $R_{yx_i.\mathbf{X}_k}^2$ for each i ($i = 1, 2, \dots, p$) gives the above decomposition. This averaging process over all permutations is carried out in the absence of a "natural" ordering of the regressors which would yield the $R_{yx_i.\mathbf{X}_k}^2$ of that ordering as measures of their relative importance (Kruskal (1987)). This decomposition involves a large amount of calculations even for moderately large p 's.

In the next Section our proposal for decomposition in the parameter space using a statistical approach will be presented and in Section 4 a description of the sample R^2 decomposition will be given from a geometrical point of view, followed by discussion and a simulation study comparing some measures of relative importance of regressors.

3. The Proposed Decomposition

Let us assume temporarily, that x_1, x_2, \dots, x_p all have zero expectation, unit variance and known correlation matrix Σ , $\text{cor}(x_i, x_j) = \text{cov}(x_i, x_j) = \Sigma_{ij}$. Assume also that $E(y) = 0$, $\text{var}(y) = 1$ and the correlations $\text{cor}(y, x_i) = \rho_{yx_i}$, $i = 1, 2, \dots, p$, are also known. It should be noted that the problem of decomposing $R^2 = \rho'_{yx} \Sigma^{-1} \rho_{yx}$, where $\rho_{yx} = (\rho_{yx_1}, \dots, \rho_{yx_p})'$, into components associated with the regressors, exists even in the present set up of known ρ_{yx} and Σ .

We are given the regression equation $Ey = \Sigma_i \beta_i x_i$. Let z_1, z_2, \dots, z_p be the (unique) linear transformation of x_1, x_2, \dots, x_p which minimizes $\Psi = E \Sigma_i (z_i - x_i)^2$ subject to the condition that the z_i be uncorrelated with expectation 0 and variance 1. This is a transformation to a new uncorrelated set of regressors

z_i , which are nearest to the corresponding x_i in the sense that the expected sum of squared deviations Ψ between the pairs (x_i, z_i) is minimized. Let the transformation be written in the form $x_i = \sum_j a_{ij}z_j$, $i = 1, 2, \dots, p$, and the regression in the new variables as

$$Ey = \sum_j c_j z_j, \quad \text{with} \quad c_j = \sum_i \beta_i a_{ij}, \quad (2)$$

then by the definition of z_1, z_2, \dots, z_p it follows that

$$a_{ij} = \rho_{z_i x_j} \quad \text{and} \quad c_j = \rho_{y z_j}. \quad (3)$$

Our proposal assigns to x_i , in the decomposition of R^2 , the component

$$v_i = \sum_j a_{ij}^2 c_j^2. \quad (4)$$

v_i is preferred over taking c_i^2 itself as the component assigned to x_i , because although each z_j is closest to the corresponding x_j in the above sense, it is still correlated in general to the other x_j 's with the coefficients a_{ij} ; therefore a component of a reasonable decomposition should be composed of a linear combination of all the c_j^2 s with coefficients which reflect the magnitude of these correlations. Further aspects of this proposal will be discussed in later sections.

It remains to obtain the transformation matrix $A = [a_{ij}]$. Since $\mathbf{x} = \mathbf{A}\mathbf{z}$, where $\mathbf{x} = (x_1, x_2, \dots, x_p)'$ and $\mathbf{z} = (z_1, z_2, \dots, z_p)'$, we have $E\mathbf{x}\mathbf{x}' = \mathbf{A}\mathbf{A}' = \Sigma$, the correlation matrix of the x 's. (Note that \mathbf{x} is a random $p \times 1$ vector, whereas the x_i 's are $n \times 1$ vectors of observations on the variable x_i , $i = 1, \dots, p$.) We have to choose \mathbf{A} among the matrix square roots of Σ in a way that it should minimize Ψ . It will now be shown that \mathbf{A} must be the symmetrical square root of Σ , i.e.

$$\mathbf{A} = \mathbf{Q}'\mathbf{D}_\alpha\mathbf{Q}, \quad (5)$$

where $\Sigma = \mathbf{Q}'\mathbf{D}_\alpha^2\mathbf{Q}$ is the principal decomposition of Σ , \mathbf{Q} is an orthogonal matrix whose rows are the eigenvectors of Σ and \mathbf{D}_α^2 is diagonal with the eigenvalues $\alpha_1^2, \alpha_2^2, \dots, \alpha_p^2$ of Σ along its diagonal.

The general form of a matrix square root of Σ is $\mathbf{A} = \mathbf{Q}'\mathbf{D}_\alpha\mathbf{G}$, with an arbitrary orthogonal matrix \mathbf{G} . We may express Ψ as

$$\begin{aligned} \Psi &= E(\mathbf{z} - \mathbf{x})'(\mathbf{z} - \mathbf{x}) = \mathbf{I} + \Sigma - 2E\mathbf{z}'\mathbf{x} = \mathbf{I} + \Sigma - 2E\mathbf{z}'\mathbf{A}\mathbf{z} \\ &= \mathbf{I} + \Sigma - 2\text{tr}E\mathbf{z}\mathbf{z}'\mathbf{A} = \mathbf{I} + \Sigma - 2\text{tr}\mathbf{A}. \end{aligned}$$

Thus, to minimize Ψ we have to maximize $\text{tr}\mathbf{A} = \text{tr}\mathbf{Q}'\mathbf{D}_\alpha\mathbf{G} = \text{tr}\mathbf{D}_\alpha\mathbf{G}\mathbf{Q}'$. To prove that $\mathbf{G} = \mathbf{Q}$ maximizes Ψ we state the following lemma.

Lemma. *If D is a diagonal matrix with positive diagonal elements, the trace of (DH) is maximized over the set of all orthogonal matrices $\{H\}$ when $H = I$ the identity matrix.*

This follows easily from $\text{tr}(DH) = \sum_i d_{ii} h_{ii}$ and $h_{ii} \leq 1$ for orthogonal H , with equality signs occurring for all i iff $H = I$.

The lemma implies that $\text{tr} D_\alpha GQ'$ with all $\alpha_i > 0$, is maximized when the orthogonal matrix $GQ' = I$, i.e., when $G = Q$. It will be required without loss of generality that $\alpha_i > 0$ ($1 \leq i \leq p$), because otherwise, to maximize $\text{tr} A$, all rows of G and Q corresponding to negative α_i must satisfy $g'_i q_i = -1$, or $g_i = -q_i$, resulting in the same A as given by (5).

Thus, A is symmetric and $A^2 = \Sigma$, hence $\sum_k a_{ik}^2 = \sum_k a_{kj}^2 = 1$ for $1 \leq i \leq p$ and $1 \leq j \leq p$, because these sums are equal to the diagonal elements of the correlation matrix Σ . It follows that the sum of components

$$\sum_i v_i = \sum_i \sum_k a_{ik}^2 c_k^2 = \sum_k \left(\sum_i a_{ik}^2 \right) c_k^2 = \sum_k c_k^2 = \sum_k \rho_{yz_k}^2 = R^2, \quad (6)$$

because the presentation of the regression as a regression on the uncorrelated set z_1, z_2, \dots, z_p gives an exact decomposition of R^2 .

To prove that decomposition (4) satisfies the OC criterion, suppose the x consists of, say, m uncorrelated subsets of variables. By rearranging x , Σ will become a block diagonal matrix, and A and Σ^{-1} will be also block diagonal with blocks at the same positions. Let the subsets be given by V_1, V_2, \dots, V_m , where $i \in V_k$ if x_i belongs to the k th subset. Thus, $R^2 = \rho'_{yx} \Sigma^{-1} \rho_{yx}$, will decompose into these subsets, say, $R^2 = R_1^2 + R_2^2 + \dots + R_m^2$, where R_k^2 is constructed only by ρ_{yx_j} 's with $j \in V_k$. Also from (2) and (4) it can be seen, that the z_i and the components v_i with $j \in V_k$, will also be formed only from a_{ji} with both indices belonging to V_k . Hence, by the same reasoning as used above for R^2 itself, it follows that

$$\sum_{i \in V_k} v_i = R_k^2. \quad (7)$$

It should be noted that (6) and (7) were proved using the special form (5) of A , and they will not hold in general for any square root of Σ . The principal decomposition of Σ is essentially unique when all eigenvalues are different. However, even when not all α_i are different, A as given by (5) and A^{-1} are always unique.

The decomposition (4) can be written using the relation $z = A^{-1}x$ as

$$v_i = \sum_j \left(\sum_k \rho_{yx_k} a^{kj} \right)^2 a_{ij}^2, \quad (8)$$

where a^{kj} is the kj th element of A^{-1} . This form involves only the correlations ρ_{yx_k} and $a_{ij} = \rho_{x_i x_j}$.

In the following, the general case will be considered when the expectation and variance of the variables are unknown. In this setup, the component assigned to x_i in the decomposition of R^2 will be estimated by (8) or (4), where the matrix $\mathbf{A} = [a_{ij}]$ now stands for the symmetric square root of the matrix of sample correlations $r_{x_i x_j}$ and $r_{y x_k}$ replaces $\rho_{y x_k}$.

4. The Sample R^2 Decomposition and its Geometrical Representation

Given a sample of $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p)$, and \mathbf{y} , we assume, without loss of generality, that they are standardized with zero means and unit lengths. In analogue to the previous procedure, we construct an orthonormal basis $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_p$ for the p dimensional subspace $V(\mathbf{X})$ in R^n spanned by $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$, with the following properties: (a) $\sum_i \mathbf{x}'_i \mathbf{z}_i$ is maximized; (b) $\mathbf{x}'_i \mathbf{z}_j = \mathbf{x}'_j \mathbf{z}_i$ for all $i, j = 1, 2, \dots, p$; and (c) $\delta = \|\mathbf{X}'\mathbf{X} - \mathbf{Z}'\mathbf{X}\|$ is minimized, where $\mathbf{Z} = (\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_p)$. Property (a) states that \mathbf{z}_i 's are closest to their corresponding \mathbf{x}_i 's in the sense that the sum of $\cos(\mathbf{x}_i, \mathbf{z}_i)$ is maximized. Property (b) shows that the \mathbf{z}_i 's are positioned around the \mathbf{x}_i 's in a symmetrical configuration so that for any pair i, j the angle between \mathbf{x}_i and \mathbf{z}_j is equal to that between \mathbf{x}_j and \mathbf{z}_i . Finally, property (c) states that $\mathbf{z}'_i \mathbf{x}_j$ is closest to $\mathbf{x}'_i \mathbf{x}_j$ in the sense that their sum of squares of deviation over all $i, j = 1, 2, \dots, p$ is minimized.

Let \mathbf{X} be presented by the basis in the form $\mathbf{X} = \mathbf{Z}\mathbf{A}$; then, analogous to the previous derivation and using the lemma, it turns out from any of the requirements (a) or (c), that \mathbf{A} is the symmetric matrix square root of $\mathbf{X}'\mathbf{X} = \mathbf{R}_{XX}$, the matrix of sample correlations between the x 's. Requirement (b) is clearly satisfied by using this matrix \mathbf{A} . In the sample all eigenvalues of $\mathbf{X}'\mathbf{X}$ are different with probability 1.

The decomposition of each \mathbf{x}_i along the directions of the \mathbf{z}_j 's is $\mathbf{x}_i = \sum_j \mathbf{z}_j \mathbf{z}'_j \mathbf{x}_i$. Similarly, each correlation coefficient $r_{y x_i} = \mathbf{y}' \mathbf{x}_i$ can be written in the form

$$r_{y x_i} = \sum_j \mathbf{y}' \mathbf{z}_j \mathbf{z}'_j \mathbf{x}_i = \sum_j c_j a_{ji}, \tag{9}$$

where $c_j = \mathbf{y}' \mathbf{z}_j = r_{y z_j}$.

In view of properties (a), (b) and (c) of the set of vectors \mathbf{z}_j , we may consider $\mathbf{y}' \mathbf{z}_i \mathbf{z}'_i \mathbf{x}_i = c_i a_{ii}$ as the "specific" component of \mathbf{x}_i , and $\mathbf{y}' \mathbf{z}_j \mathbf{z}'_j \mathbf{x}_i = c_j a_{ji}$ as the "common" component of \mathbf{x}_i and \mathbf{x}_j ($j \neq i = 1, 2, \dots, p$) in the correlation coefficient between \mathbf{y} and \mathbf{x}_i .

Likewise, we define the component assigned to \mathbf{x}_i in the decomposition of $R^2_{y(\mathbf{X})}$ as

$$r^2_{y x_i}(\mathbf{X}) = \sum_j (\mathbf{y}' \mathbf{z}_j \mathbf{z}'_j \mathbf{x}_i)^2 = \sum_j (c_j a_{ji})^2. \tag{10}$$

Here we changed the component's notation to indicate its dependence on the whole set of regressors. It contains, as above, a specific component due to \mathbf{x}_i and common components with all other \mathbf{x}'_j $j \neq i$.

Using $\mathbf{Z} = \mathbf{X}\mathbf{A}^{-1}$, we can express $\mathbf{c} = (c_1, \dots, c_p)$, too, as functions of correlations only, in the form $\mathbf{c} = \mathbf{Z}'\mathbf{y} = \mathbf{A}^{-1}\mathbf{X}'\mathbf{y} = \mathbf{A}^{-1}\mathbf{r}_{\mathbf{X}\mathbf{y}}$; hence $r_{y\mathbf{x}_i}^2(\mathbf{X})$ can be expressed in the form

$$r_{y\mathbf{x}_i}^2(\mathbf{X}) = \sum_j \left(\sum_k r_{y\mathbf{x}_k} a^{kj} \right)^2 a_{ji}^2, \quad (11)$$

where a^{kj} is the kj th element of \mathbf{A}^{-1} .

Proofs that this is a decomposition, i.e., $\sum_i r_{y\mathbf{x}_i}^2(\mathbf{X}) = R^2$, and that it satisfies Criterion OC when component sets of the sample vectors are mutually orthogonal are completely analogous to those given above for the population, and will not be repeated here.

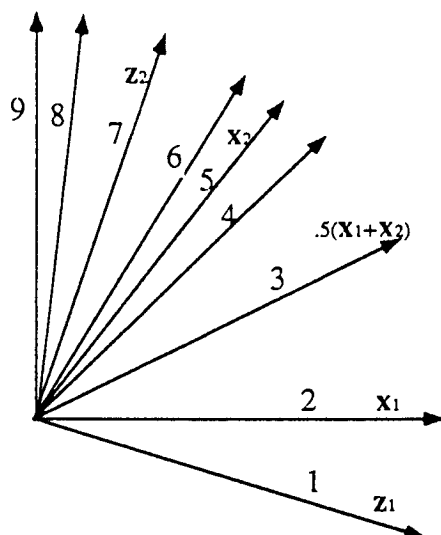
5. Discussion

(a) Let $\mathbf{R}_{\mathbf{X}\mathbf{X}} \rightarrow \mathbf{I}$ in the form $c\mathbf{R}_{\mathbf{X}\mathbf{X}} + (1-c)\mathbf{I}$, where $c \rightarrow 0$. This is equivalent to having $\mathbf{D}_\alpha^2 \rightarrow \mathbf{I}$ in the form $c\mathbf{D}_\alpha^2 + (1-c)\mathbf{I}$ while using the *same* matrix \mathbf{Q} , since $\mathbf{Q}'\{c\mathbf{D}_\alpha^2 + (1-c)\mathbf{I}\}\mathbf{Q} = c\mathbf{R}_{\mathbf{X}\mathbf{X}} + (1-c)\mathbf{I}$. Therefore, if we also leave \mathbf{P} unchanged, the coordinate system $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_p$ may remain, by this method, preserved, while each vector \mathbf{x}_i will tend to its corresponding coordinate vector, and $r_{y\mathbf{x}_i}^2(\mathbf{X})$ will converge to $(\mathbf{y}'\mathbf{z}_i)^2 = c_i^2 = (\sum_k r_{y\mathbf{x}_k} a^{ki})^2$. Thus, the special coordinate system has a further property, that it is the set to which the vectors \mathbf{x}_i can smoothly tend when uncorrelatedness is approached in the simple way of reducing simultaneously all correlation coefficients towards zero by the same proportion.

(b) A few special cases will now be considered to illustrate the proposed decomposition.

1. $p = 2$. In this case, to construct orthogonal \mathbf{z}_1 and \mathbf{z}_2 satisfying $\mathbf{z}'_1\mathbf{x}_2 = \mathbf{z}'_2\mathbf{x}_1$, the angle between \mathbf{x}_1 and \mathbf{x}_2 has to be enlarged (or constricted) symmetrically with respect to their bisector (Fig.1) to an orthogonal one. Figure 1 also shows the changes occurring in the components of two measures when the position of $\hat{\mathbf{y}}$ is changed in relation to a specific configuration of \mathbf{x}_1 and \mathbf{x}_2 . The cases $\hat{\mathbf{y}} = \mathbf{x}_1$ and $\hat{\mathbf{y}} = \mathbf{x}_2$ are unrealistic because in such cases the model will consist of one regressor only. They can be viewed as limiting cases $\hat{\mathbf{y}} \rightarrow \mathbf{x}_1$ or $\hat{\mathbf{y}} \rightarrow \mathbf{x}_2$ when R^2 also tends to 1 and therefore even a small contribution of the other regressor is important. The case when one of the simple correlations is zero, say $r_{y\mathbf{x}_1}^2 = 0$ ($\hat{\mathbf{y}}$ in position 9 in Figure 1), is instructive. Usually one will have, even in this case, $R^2 > r_{y\mathbf{x}_2}^2$. Therefore, it is justified to state that \mathbf{x}_1 has also contributed to R^2 and thus to have a positive component for \mathbf{x}_1 . Indeed, \mathbf{x}_1 is entitled to one

half of the difference $R^2 - r_{yx_2}^2$. In a similar way the rationale behind the mode of changes in $r_{yx_i}^2(\mathbf{X})$ can be seen to be reasonable.



position of \hat{y}	$r_{yx_1}^2$	$r_{yx_2}^2$	$r_{yx_1}^2(\mathbf{X})$	$r_{yx_2}^2(\mathbf{X})$	$\hat{y}'\hat{y}_1$	$\hat{y}'\hat{y}_2$
1	.90	.10	.90	.10	1.125	-.125
2	1.00	.36	.82	.18	1.00	.00
3	.80	.80	.50	.50	.50	.50
4	.50	.98	.26	.74	.125	.875
5	.36	1.00	.18	.82	.00	1.00
6	.26	.99	.14	.86	-.07	1.07
7	.10	.90	.10	.90	-.125	1.125
8	.01	.73	.14	.86	-.07	1.07
9	.00	.64	.18	.82	.00	1.00

Figure 1. Squared correlation coefficients and components of R^2 (in units of R^2) according to two definitions, $r_{yx_i}^2(\mathbf{X}) = \bar{r}_i^2$ and $\hat{y}'\hat{y}_i$, for different directions of \hat{y} in a specific case with $p = 2$ and $\cos(x_1, x_2) = 0.6$ (deviation measure from orthogonality $\Delta = 0.7$).

In this case $r_{yx_i}^2(\mathbf{X})$ ($i = 1, 2$) can be expressed by simpler terms in the form $r_{yx_i}^2(\mathbf{X}) = (r_{yx_i}^2 + R_{yx_i, x_j}^2)/2$, $i \neq j = 1, 2$. It can be derived in the following way. Due to $a_{12} = a_{21}$ and the positive definiteness of \mathbf{A} , we also have $a_{11} = a_{22}$ because $a_{11}^2 + a_{12}^2 = a_{21}^2 + a_{22}^2 = 1$. Also by (2) and (3) the differences between the components of R^2 and the simple squared correlations are equal i.e. $r_{yx_1}^2(\mathbf{X}) - r_{yx_1}^2 = 2c_1c_2a_{11}a_{12} = r_{yx_2}^2(\mathbf{X}) - r_{yx_2}^2 = 2c_1c_2a_{21}a_{22}$. Therefore, $r_{yx_1}^2 + R_{yx_1, x_2}^2 = r_{yx_1}^2 + R^2 - r_{yx_2}^2 = r_{yx_1}^2(\mathbf{X}) + R^2 - r_{yx_2}^2(\mathbf{X}) = 2r_{yx_1}^2(\mathbf{X})$ and the result follows.

2. \hat{y} forms equal angles with all z_i 's, or $c_1 = c_2 = \dots = c_p$ ($= c_0$, say). In

this case $\hat{\mathbf{y}}$ is in the center of the set $\{\mathbf{z}_i\}$ (i.e., $\hat{\mathbf{y}} = R/\sqrt{p}\mathbf{Z}\mathbf{1}$) and all $r_{y x_i}^2(\mathbf{X})$ are equal because $r_{y x_i}^2(\mathbf{X}) = \sum_j c_j^2 a_{ij}^2 = c_0^2 \sum_j a_{ij}^2 = c_0^2 = R^2/p$. Such equality does not necessarily hold for the $r_{y x_i}^2$'s.

3. All $r_{x_i x_j}$'s are equal to (say) r so that the \mathbf{x}_i 's form a symmetric pencil of vectors with axis of symmetry $\mathbf{X}\mathbf{1}_p$. In this case \mathbf{z}_j ($j = 1, 2, \dots, p$) will form equal angles with all \mathbf{x}_i , $i \neq j$, and have the same axis of symmetry. Also, it is easily shown that $\mathbf{A} = \sqrt{(1-r)}\{\mathbf{I}_p + p^{-1}(q-1)\mathbf{J}_p\}$, where $q = \sqrt{(1+pr-r)}/\sqrt{(1-r)}$, and \mathbf{J}_p is a $p \times p$ matrix of 1's. Therefore, in this case, $r_{y x_i}^2(\mathbf{X})$ can also be explicitly expressed in the form $r_{y x_i}^2(\mathbf{X}) = (1-r)(q-1)^2 R^2/p^2 + \{1 + 2(q-1)/p\}\{\mathbf{y}'\mathbf{x}_i - (1-q^{-1})\sum_j \mathbf{y}'\mathbf{x}_j/p\}^2$, $i = 1, 2, \dots, p$.

(c) Pratt (1987) discussed a symmetrical situation, which is equivalent to assuming both 2. and 3. above. The components of two subsets consisting of m and n regressors, respectively, are in the ratio of $m : n$, as follows from 2., but in a regression on two variables which are the sums of the m and n regressors, respectively, our measure yields components in a ratio which is a monotonic increasing function of m/n , whereas Pratt's measure $\hat{\mathbf{y}}'\hat{\mathbf{y}}_i$ maintains the ratio $m : n$. Four out of Pratt's six criteria are also satisfied by our measure. The exceptions are his postulate (A3), which requires the above mentioned ratio $m : n$, and (A6), which requires invariance of subsets' components under nonsingular linear transformation of any one of the other subsets. Neither are these postulates satisfied by \bar{r}_i^2 of (1). It can be argued from the geometrical point of view, that this assumption is not so "natural" if the subsets are not mutually orthogonal, because change in position of the individual vectors of the transformed subset in relation to other subsets, changes the overall configuration, which is expressed in the coordinate set \mathbf{Z} and hence, in our opinion, in the whole decomposition it implies.

(d) Maximum likelihood estimates of $\rho_{y x_j}$ ($i = 1, 2, \dots, p$) under a multivariate normal model in which the \mathbf{x}_i 's are assumed to be uncorrelated, can also be proposed as candidates (after squaring) for components in the decomposition of R^2 . They can be derived in the form $\hat{\rho}_{y x_j} = \sum_j q(\mathbf{X}'\mathbf{X})^{ij} \mathbf{y}'\mathbf{x}_j = \sum_j q \mathbf{A}^{ij} c_j$, where q is a factor of proportionality and $(\mathbf{X}'\mathbf{X})^{ij}$ is the ij th element of $(\mathbf{X}'\mathbf{X})^{-1}$. Since the i th row of \mathbf{A}^{-1} is orthogonal to all columns of \mathbf{A} except for the i th, these estimates, after squaring, are proportional to $R_{y x_i, \mathbf{X}_{-i}}^2$, which were discussed in Section 2 and rejected from serving as components in the decomposition of R^2 . In fact, our aim is to find a decomposition which would assign to each variable its 'fair' portion of R^2 in the presence of nonzero correlations, instead of assuming that they are zero.

(e) The matrices \mathbf{Q} and \mathbf{D}_α occur also in principal components regression (e.g. Massy (1965) and Mansfield et al. (1977)), which is done to facilitate and simplify calculations in cases of collinearity, and to base upon it considerations about

reducing the number of regressors.

6. Decomposition in ANOVA Problems

If we relax the requirement of linear independence between $\mathbf{x}_1, \mathbf{x}_1, \dots, \mathbf{x}_p$ then definition (11) of $r_{y\mathbf{x}_i}^2(\mathbf{X})$ can still hold, with a^{kj} being the kj th element of the generalized inverse of \mathbf{A} obtained by inverting the positive elements of \mathbf{D}_α and leaving the zero elements unchanged (a Moore Penrose type generalized inverse). This may be advantageous when there is, *a priori*, a “natural” set of regressors into which the variance of \mathbf{y} is to be decomposed despite their linear interdependence.

It is well known (cf. Rodgers and Nicewander (1988)) that, in a one-way ANOVA setting, the squared multiple correlation coefficient between y and the columns of the design matrix \mathbf{X} is a monotone function of the F statistic testing the hypothesis of equality of the population means. This is an example of the case where the standardized \mathbf{x}_i 's are linearly dependent. It is interesting to check the meaning of the individual components of R^2 in terms of contrasts between the populations' means. Let the number of populations be k and the sample size in each population be n . We have $R_{y\mathbf{x}_i, \mathbf{X}_{-i}}^2 = 0, i = 1, 2, \dots, k$, due to the linear dependence, but the three measures (i), (iii) and (iv) from Section 2, namely $r_{y\mathbf{x}_i}^2, \hat{\mathbf{y}}_i' \hat{\mathbf{y}}_i$ and $\hat{\mathbf{y}}' \hat{\mathbf{y}}_i$ are proportional to the same squared standardized contrast: $b_i^2/S_{yy} = n(y_i - y_{..})^2/\Sigma \Sigma (y_{ij} - y_{..})^2$, with $\Sigma b_i^2/S_{yy} = R^2$, using the customary dot notation. This happens because of the special form of $\mathbf{X}'\mathbf{X} = \mathbf{I}_k - \frac{1}{k}\mathbf{J}_k$. For $r_{y\mathbf{x}_i}^2(\mathbf{X})$ we get $r_{y\mathbf{x}_i}^2(\mathbf{X}) = [R^2/k + (k - 2)b_i^2/S_{yy}]/(k - 1)$, which is a weighted average of R^2/k and b_i^2/S_{yy} . Note that for $k = 2$ the four measures coincide.

In the situation when y_i coincides with the mean of the other samples, $y_i = y_{..}$, the three measures associating component zero to \mathbf{x}_i , give the intuitively correct value, because \mathbf{x}_i does not contribute to R^2 at all. On the other hand, when y_i is merely close to the mean of the other samples, then the contribution of population i to R^2 should also reflect the increase in R^2 obtained from $\sum_{j \neq i} b_j^2/S_{yy}$ due to its effect on $y_{..}$. Thus, in general $r_{y\mathbf{x}_i}^2(\mathbf{X})$ is more justified than the other measures in the ANOVA context.

7. A Measure of Deviation from Orthogonality

Given \mathbf{X} , we found \mathbf{Z} that minimizes $\delta = \|\mathbf{X}'\mathbf{X} - \mathbf{Z}'\mathbf{X}\|$, over all choices of orthonormal matrices \mathbf{Z} , which is : $\delta = \|\mathbf{D}_\alpha^2 - \mathbf{D}_\alpha\| = \sqrt{\Sigma_i (\alpha_i^2 - \alpha_i)^2} = \sqrt{\Sigma_i \alpha_i^2 (\alpha_i - 1)^2}$. When the columns of \mathbf{X} are orthogonal, $\mathbf{X}'\mathbf{X} = \mathbf{I}$ and all $\alpha_i (= +\sqrt{\alpha_i^2})$ are 1, hence $\mathbf{Z} = \mathbf{X}$ and $\delta = 0$. Since $\Sigma \alpha_i^2 = \text{tr}(\mathbf{X}'\mathbf{X}) = p$, the supremum of δ over different \mathbf{X} matrices occurs when $\max_i \alpha_i^2 \rightarrow p$ and all

other α_j 's $\rightarrow 0$. Thus, we can normalize δ by putting $\Delta = \delta/(p - \sqrt{p})$, yielding $0 \leq \Delta < 1$. The deviation of Δ from zero towards 1 can be used as a measure of deviation of the columns of \mathbf{X} from orthogonality. Δ is a function of the characteristic roots α_i^2 of the matrix of correlations; therefore it is a meaningful measure in both the sample and the parameter space. Since the decompositions mentioned in Section 2 all agree when \mathbf{X} is orthogonal it is interesting to find out in a quantitative way their trend to disagree as Δ approaches 1. This is done in the next Section.

8. Comparisons of Measures: Simulations

Numerical calculations of the various measures on real data sets are of limited value, for the decision as to which of the measures assesses, fairly to the regressors, relative importances in a sense which is relevant to the data, is a matter of subjective judgment.

A small Monte Carlo simulation was carried out to calculate correlation coefficients between components of six measures of relative importance for $p = 2, 3$ and 4, to investigate their closeness under different conditions. 2000 random $\mathbf{X}(p+1) \times p$ matrices and random $p \times 1 \mathbf{c} = \mathbf{z}'\mathbf{y}$ vectors were generated, all i.i.d. $U[0, 1]$. In each case the measure of deviation from orthogonality, Δ , and the components for six measures were calculated.

Table 1 shows the number (n) of cases that fall in each Δ range and the average of correlation coefficients calculated for each sample between the components $r_{yx_i}^2(\mathbf{X})$ and each of the following measures: SIMPLE = $r_{yx_i}^2$, MARGINAL = $R_{yx_i, \mathbf{X}_{-i}}^2$, STANDARD = $\hat{\mathbf{y}}_i'\hat{\mathbf{y}}_i$, AVERAGE = \bar{r}_i^2 from (1), PRATT = $\hat{\mathbf{y}}_i'\hat{\mathbf{y}}_i$ and NNPRATT, which is the same as PRATT, but includes only the cases when all the components are non-negative. The standard deviations are also presented. For $p = 2$, all correlation coefficients were 1, indicating a complete agreement between all measures concerning the order of importance of the two regressors. Therefore, it was omitted from the table. However, the average mean square of deviations between $r_{yx_i}^2(\mathbf{X})$ and the other measures changed, for $p = 2$, with increasing Δ from .0004 to .04, with small differences among the six measures. For $p > 2$, even the same ordering with different relative components results in $r < 1$, and the value of r reflects the magnitude of agreement among the measures.

It is evident that as $\Delta \rightarrow 1$ there is a sharp decrease in the correlation coefficients. They also decrease somewhat with increasing p , especially for small Δ values. Among the measures, AVERAGE is closest to $r_{yx_i}^2(\mathbf{X})$ followed by NNPRATT, PRATT, SIMPLE, MARGINAL and STANDARD in decreasing order. Note that n_{NN} , the number of cases for NNPRATT, decreased sharply with increasing p and Δ , indicating that some negative components for PRATT occur quite frequently. The table shows that for small p and Δ , several measures

Table 1. Average correlation coefficients (\pm standard deviations) between the components of the proposed measure of relative importance and those of six other measures from a Monte Carlo study with $p = 3$ and 4 and at different ranges of the nonorthogonality measure Δ . The number of cases (n) at each Δ range is also presented.

$\Delta =$.0-0.3	.3-0.4	.4-0.5	.5-0.6	.6-0.7	.7-0.8	.8-0.9	.9-1.0
<u>measure*</u>								
	$p = 3$							
SIMPLE	.94 \pm .22	.70 \pm .53	.67 \pm .58	.67 \pm .56	.70 \pm .54	.62 \pm .62	.74 \pm .50	.18 \pm .67
MARGINAL	.97 \pm .10	.80 \pm .40	.64 \pm .61	.63 \pm .60	.61 \pm .60	.52 \pm .70	.49 \pm .74	.66 \pm .60
STANDARD	.97 \pm .07	.70 \pm .48	.46 \pm .70	.54 \pm .62	.45 \pm .64	.34 \pm .73	.28 \pm .73	.44 \pm .63
AVERAGE	1.00 \pm .01	.85 \pm .41	.84 \pm .44	.82 \pm .44	.83 \pm .43	.80 \pm .49	.88 \pm .32	.83 \pm .31
PRATT	1.00 \pm .01	.88 \pm .28	.73 \pm .41	.74 \pm .34	.71 \pm .32	.64 \pm .41	.57 \pm .43	.53 \pm .46
NNPRATT	1.00 \pm .01	.90 \pm .27	.81 \pm .40	.77 \pm .33	.79 \pm .33	.56 \pm .57	.47 \pm .64	.32 \pm .55
$n(n_{NN})$	125(100)	456(174)	554(160)	345(85)	227(38)	146(32)	107(13)	40(8)
	$p = 4$							
SIMPLE	.77 \pm .36	.71 \pm .41	.67 \pm .44	.65 \pm .48	.65 \pm .50	.63 \pm .46	.70 \pm .47	
MARGINAL	.73 \pm .44	.62 \pm .51	.58 \pm .52	.58 \pm .52	.54 \pm .53	.45 \pm .62	.16 \pm .52	
STANDARD	.60 \pm .49	.49 \pm .54	.44 \pm .54	.44 \pm .51	.39 \pm .49	.41 \pm .52	.06 \pm .69	
AVERAGE	.93 \pm .19	.91 \pm .22	.90 \pm .24	.91 \pm .20	.91 \pm .15	.89 \pm .23	.96 \pm .05	
PRATT	.82 \pm .31	.73 \pm .33	.70 \pm .31	.65 \pm .32	.60 \pm .33	.66 \pm .30	.45 \pm .35	
NNPRATT	.90 \pm .25	.83 \pm .29	.80 \pm .30	.70 \pm .37	.66 \pm .45	.87 \pm .08	-	
$n(n_{NN})$	316(92)	673(120)	607(77)	282(29)	97(9)	24(2)	7(0)	

*SIMPLE = $r_{yx_i}^2$, MARGINAL = $R_{yx_i \cdot X_{-i}}^2$, STANDARD = $\hat{y}'_i \hat{y}_i$, AVERAGE = \bar{r}_i^2
 PRATT = $\hat{y}' \hat{y}_i$; NNPRATT is the same as PRATT, but discards all cases when any component is negative; n_{NN} is the number of cases for NNPRATT.

give essentially the same decomposition and any of them may be used to assess the relative importance of the regressors. For larger p or Δ the measures vary considerably, and one should use a measure only after its relevance has been established.

In the example of Yule's data, for which Pratt (1987) calculated the components of his measure, there are $p = 3$ regressors and $\Delta = 0.34$. Hence, we found that, in accordance with Table 1, all the above measures are highly correlated and gave similar estimates of relative importances.

Acknowledgements

The author is very grateful to the referees for their many helpful comments and suggestions. This research was partly sponsored by the US-Israel Binational Agricultural Research and Development Fund, Project No. US-334-81. Partial

assistance was obtained from Natural Sciences and Engineering Research Council of Canada.

References

- Kruskal, W. (1984). Relative importance of determiners. *Questio* 8, 39-45.
- Kruskal, W. (1987). Relative importance by averaging over orderings. *Amer. Statist.* 41, 6-10.
- Kruskal, W. and Majors, R. (1989). Concepts of relative importance in recent scientific literature. *Amer. Statist.* 43, 2-6.
- Mansfield, E. R., Webster, J. T. and Gunst, R. F. (1977). An analytic variable selection technique for principal component regression. *Appl. Statist.* 26, 34-40.
- Massy, W. F. (1965). Principal components regression in exploratory statistical research. *J. Amer. Statist. Assoc.* 60, 234-256.
- Pratt, J. W. (1987). Dividing the indivisible: using simple symmetry to partition variance explained. In *Proceedings of Second International Tampere Conference in Statistics* (Edited by T. Pukkila and S. Puntanen), 245-260. University of Tampere, Tampere, Finland.
- Rodgers, J. L. and Nicewander, W. A. (1988). Thirteen ways to look at the correlation coefficient. *Amer. Statist.* 42, 59-66.
- Williams, E. J. (1978). Postscripts to "Linear hypothesis: regression". In *International Encyclopedia of Statistics* (Edited by W. H. Kruskal and J. M. Tanur), 537-541. Free Press, New York, NY.

Department of Statistics and Operations Research, Agricultural Research Organization, Bet Dagan 50250, P.O. Box 6, Israel.

(Received June 1991; accepted December 1992)