

STICHTING
MATHEMATISCH CENTRUM
2e BOERHAAVESTRAAT 49
AMSTERDAM
AFDELING MATHEMATISCHE STATISTIEK

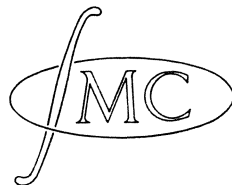
Report S 319 (VP 23)

On the selection of independent variables in a regression equation

Preliminary report

by

J. Oosterhoff



December 1963

The Mathematical Centre at Amsterdam, founded the 11th of February, 1946, is a non-profit institution aiming at the promotion of pure mathematics and its applications, and is sponsored by the Netherlands Government through the Netherlands Organization for Pure Research (Z.W.O.) and the Central National Council for Applied Scientific Research in the Netherlands (T.N.O.), by the Municipality of Amsterdam and by several industries.

1. Introduction

In many applications of regression theory it may be of interest to reduce the set of independent variables to a smaller subset. If we have m independent variables and we want to retain only k variables in the regression equation ($k < m$), the problem arises how to select these k variables. Here we consider different powers of a mathematical variable as different independent variables. For notational convenience we define a k -subset as a subset of k independent variables. A k -subset will be called better than another k -subset, if the sum of squares due to regression on the first subset is larger than the sum of squares due to regression on the second subset. A k -subset is an optimal k -subset, if no better k -subset exists. One would like to have available a computationally easy method by which optimal k -subsets are unfailingly designated, but unfortunately such a method is as yet not known. Of course the sums of squares due to regression on all possible $\binom{m}{k}$ k -subsets may be computed, but this is by no means a quick method for large $\binom{m}{k}$. Computational labour is still more formidable, because we usually compute the sums of squares due to regression on k -subsets for different k , and let k depend on the outcome.

In an expository paper [1] H.C. HAMAKER describes two techniques, called forward selection and backward elimination, which often lead to optimal k -subsets. We give a brief outline of both methods. To fix ideas, introduce $m+1$ vectors X_1, X_2, \dots, X_m and \underline{Y} of dimension n , where X_i denotes the vector of the n (non-random) observations on the i^{th} independent variable x_i ($i=1, \dots, m$), and \underline{Y} the vector of the corresponding n observations on the dependent variable \underline{y} .

The principle of forward selection, also called stepwise regression, is as follows. At each step of the procedure we add to the independent variables already selected in the equation at former stages that independent variable, which among all remaining variables gives rise to the largest increase of the sum of squares due to regression. Continuing this process until all independent variables are selected in the equation, we successively find k -subsets to be included in the regression equation for $k=1, 2, \dots, m$.

Applying the method of backward elimination we start with the full equation containing all independent variables, and step by step remove the independent variables from our equation in the order in which the smallest decrease of the sum of squares due to regression at each step is produced, until all independent variables are eliminated from the equation. In this way we also find k -subsets for $k=m, m-1, \dots, 1$.

For computational details and a number of worked examples see [1]. Another discussion is given by P.G. MOORE in [2]. We remark that both methods are easy to apply; k -subsets are directly derived from the matrix of crossproduct sums (forward selection) or its inverse (backward elimination). Each method yields a sequence of the m independent variables, in the order in which they would be included in the equation, and the corresponding sums of squares.

It is well known that neither of these methods necessarily leads to optimal k -subsets, except in the trivial cases $k=1$ or $k=m-1$ respectively. Moreover, both methods may lead to different k -subsets. In [1] HAMAKER raised the question, whether identity of the sequences of the independent variables yielded by both methods is a sufficient condition for the optimality of the produced k -subsets for all k . We first present three fictitious numerical examples demonstrating that this condition is not sufficient. A more general example is treated in the next section.

2. Some numerical examples

In the examples (see tables I, II and III) we have four independent variables x_1, x_2, x_3, x_4 and a dependent variable y . The ten observations on each of the variables average to zero, so no constant term is needed in the regression equations. In the tables the matrices of crossproduct sums corresponding to the sets of observations are given together with the sums of squares due to regression on all subsets, the latter as fractions of the total sums of squares (rounded off to four decimal places).

TABLE I Example 1

Observations

x_1	0.71	0	0	0	0	-0.71	0	0	0	0
x_2	0.14	0.69	0	0	0	-0.14	-0.69	0	0	0
x_3	0.42	0.27	0.49	0	0	-0.42	-0.27	-0.49	0	0
x_4	0.64	-0.13	-0.07	0.27	0	-0.64	0.13	0.07	-0.27	0
y	0.57	-0.32	-0.32	0.09	0.1	-0.57	0.32	0.32	-0.09	-0.1

Matrix of Crossproduct Sums

	x_1	x_2	x_3	x_4	y
x_1	1.0082	0.1988	0.5964	0.9088	0.8094
x_2		0.9914	0.4902	-0.0002	0.5046
x_3			0.9788	0.3988	0.3002
x_4				1.0086	0.758
y					1.0158

Sums of Squares due to Regression on subsets of independent variables
in fractions of the total sum of squares

SS $\{x_1\}$ = 0.6397	SS $\{x_2, x_4\}$ = 0.8138	SS $\{x_1, x_2, x_3\}$ = 0.9644
SS $\{x_4\}$ = 0.5608	SS $\{x_1, x_2\}$ = 0.7627	SS $\{x_2, x_3, x_4\}$ = 0.9194
SS $\{x_2\}$ = 0.2528	SS $\{x_1, x_3\}$ = 0.6899	SS $\{x_1, x_2, x_4\}$ = 0.8179
SS $\{x_3\}$ = 0.0906	SS $\{x_1, x_4\}$ = 0.6439	SS $\{x_1, x_3, x_4\}$ = 0.6906
	SS $\{x_3, x_4\}$ = 0.5608	
	SS $\{x_2, x_3\}$ = 0.2563	

TABLE II Example 2

Observations

x_1	0.71	0	0	0	0	-0.71	0	0	0	0
x_2	0	0.71	0	0	0	0	-0.71	0	0	0
x_3	0.35	0.49	0.36	0	0	-0.35	-0.49	-0.36	0	0
x_4	0.28	0.21	-0.57	0.13	0	-0.28	-0.21	0.57	-0.13	0
y	0.57	0.35	-0.21	0.01	0.05	-0.57	-0.35	0.21	-0.01	-0.05

Matrix of Crossproduct Sums

	x_1	x_2	x_3	x_4	y
x_1	1.0082	0	0.497	0.3976	0.8094
x_2		1.0082	0.6958	0.2982	0.497
x_3			0.9844	-0.0086	0.5908
x_4				0.9286	0.7082
y					0.9882

Sums of Squares due to Regression on subsets of independent variables in fractions of the total sum of squares

SS $\{x_1\}$ = 0.6576	SS $\{x_3, x_4\}$ = 0.9134	SS $\{x_1, x_2, x_3\}$ = 0.9947
SS $\{x_4\}$ = 0.5466	SS $\{x_1, x_2\}$ = 0.9055	SS $\{x_1, x_2, x_4\}$ = 0.9922
SS $\{x_3\}$ = 0.3588	SS $\{x_1, x_4\}$ = 0.8560	SS $\{x_1, x_3, x_4\}$ = 0.9844
SS $\{x_2\}$ = 0.2479	SS $\{x_1, x_3\}$ = 0.7079	SS $\{x_2, x_3, x_4\}$ = 0.9713
	SS $\{x_2, x_4\}$ = 0.6272	
	SS $\{x_2, x_3\}$ = 0.3712	

TABLE III Example 3

Observations

x_1	0.71	0	0	0	0	-0.71	0	0	0	0
x_2	0.18	0.68	0	0	0	-0.18	-0.68	0	0	0
x_3	0.42	0.47	0.31	0	0	-0.42	-0.47	-0.31	0	0
x_4	0.46	-0.12	-0.39	0.36	0	-0.46	0.12	0.39	-0.36	0
y	0.57	0.29	-0.26	0.07	0.15	-0.57	-0.29	0.26	-0.07	-0.15

Matrix of Crossproduct Sums

	x_1	x_2	x_3	x_4	y
x_1	1.0082	0.2556	0.5964	0.6532	0.8094
x_2		0.9896	0.7904	0.0024	0.5996
x_3			0.9868	0.0318	0.5902
x_4				1.0154	0.708
y					1.008

Sums of Squares due to Regression on subsets of independent variables in fractions of the total sum of squares

SS $\{x_1\}$ = 0.6446	SS $\{x_2, x_4\}$ = 0.8482	SS $\{x_1, x_2, x_3\}$ = 0.9456
SS $\{x_4\}$ = 0.4897	SS $\{x_3, x_4\}$ = 0.8144	SS $\{x_1, x_2, x_4\}$ = 0.9244
SS $\{x_2\}$ = 0.3604	SS $\{x_1, x_2\}$ = 0.8115	SS $\{x_2, x_3, x_4\}$ = 0.8711
SS $\{x_3\}$ = 0.3502	SS $\{x_1, x_4\}$ = 0.7011	SS $\{x_1, x_3, x_4\}$ = 0.8164
	SS $\{x_1, x_3\}$ = 0.6641	
	SS $\{x_2, x_3\}$ = 0.3950	

Applying forward selection to the first example we successively select x_1, x_2, x_3, x_4 in the regression equation, and using backward elimination we eliminate x_4, x_3, x_2, x_1 in this order, as may be seen from the sums of squares in table I. Both methods thus lead to the sequence x_1, x_2, x_3, x_4 , and the 2-subset yielded is $\{x_1, x_2\}$. The optimal 2-subset however is $\{x_2, x_4\}$.

In the second example, application of both methods leads to the sequence x_1, x_2, x_3, x_4 , as again may be seen from the sums of squares in table II. As a 2-subset they yield $\{x_1, x_2\}$; the optimal 2-subset contains in this case the other two variables x_3, x_4 . This example has another interesting property. It might be hoped, that repeatedly substituting some variables of a k-subset for some other variables in such a way that the resulting k-subset is a better one, would lead to an optimal k-subset. This example shows that this is not true, because $\{x_1, x_2\}$ is the best 2-subset that contains either x_1 or x_2 .

In the third example both methods lead to the sequence x_1, x_2, x_3, x_4 , see table III, and yield $\{x_1, x_2\}$ as a 2-subset. Now $\{x_2, x_4\}$ and $\{x_3, x_4\}$ are both better 2-subsets.

3. A general example

We conclude with an example of a more general nature. We do not underline random variables in this section, because all considerations are in terms of fixed values of the observations.

Let k and m be fixed integers, $1 < k < m-1$, and let Y_1, Y_2, \dots, Y_{m+1} be $m+1$ orthogonal unit vectors in the n -dimensional euclidean space ($n > m+1$). Define the vectors

$$X_1 = Y_1, X_2 = Y_2, \dots, X_{k+1} = Y_{k+1},$$

$$X_j = \xi_{1j} Y_1 + \xi_{2j} Y_2 + \dots + \xi_{k+1,j} Y_{k+1} + \eta_j Y_j \quad (j=k+2, \dots, m),$$

$$Y = \omega_1 Y_1 + \omega_2 Y_2 + \dots + \omega_m Y_m + \theta Y_{m+1},$$

where $\eta, \theta, \{\xi_{ij}\}$ and $\{\omega_i\}$ are sets of real numbers satisfying

$$\sum_{i=1}^{k+1} \xi_{ij}^2 + \eta^2 = 1 \quad (j=k+2, \dots, m) \quad (1)$$

$$\sum_{i=1}^m \omega_i^2 = 1. \quad (2)$$

The vector Y is to be interpreted as the vector of observations on the dependent variable y , and the vectors X_1, X_2, \dots, X_m as the vectors of observations on the independent variables x_1, x_2, \dots, x_m .

We introduce the following condition in the above model: all vectors $X_j - \eta Y_j$ ($j=k+2, \dots, m$) are coplanar with the vectors X_{k+1} and $\omega_1 Y_1 + \omega_2 Y_2 + \dots + \omega_{k+1} Y_{k+1}$, and all these $m-k+1$ vectors have different directions. Or

$$\xi_{ij} = \gamma_j \omega_i \quad \text{and} \quad \xi_{k+1, j} = \gamma_j \omega_{k+1} + \delta_j \quad (3)$$

($i=1, \dots, k; \quad j=k+2, \dots, m$),

where $\{\gamma_j\}$ and $\{\delta_j\}$ are sets of non-zero real numbers satisfying (cf. (1))

$$\left(\frac{\gamma_j}{\delta_j}\right)^2 \sum_{i=1}^k \omega_i^2 + \left(\frac{\gamma_j}{\delta_j} \omega_{k+1} + 1\right)^2 = \frac{1-\eta^2}{\delta_j^2} \quad (4)$$

and $\gamma_s/\delta_s \neq \gamma_t/\delta_t$ for all $s \neq t$.

We now state a theorem for this model, which will permit us to draw some pertinent conclusions.

Theorem: Let $\{c_i | i=1, \dots, k\}$ and $\{d_j | j=k+2, \dots, m\}$ be given sets of different real numbers satisfying the relations

$$c_1 > c_2 > \dots > c_k > 1 \quad \text{and} \quad 0 < d_j < \left(1 + \sum_{i=1}^k c_i^2\right)^{-\frac{1}{2}} \quad \text{for all } j. \quad (5)$$

Let ϵ and ζ be arbitrarily small positive numbers ($0 < \epsilon, \zeta < 1$).

Then it is always possible to find an $\eta > 0$ and sets of coefficients $\{\omega_i | i=1, \dots, m\}$, $\{\gamma_j | j=k+2, \dots, m\}$ and $\{\delta_j | j=k+2, \dots, m\}$ such, that simultaneously (independent of the value of θ)

$$(i) \quad \frac{\omega_i}{\omega_{k+1}} = c_i \quad (i=1, \dots, k) \quad \text{and} \quad \frac{|\gamma_j|}{\delta_j} = d_j \quad (j=k+2, \dots, m),$$

(ii) forward selection and backward elimination both yield the sequence of the independent variables x_1, x_2, \dots, x_m in this order,

(iii) the sum of squares due to regression on any pair of variables from the set $\{x_l | l=k+1, \dots, m\}$ is larger than the sum of squares due to regression on the variables x_1, x_2, \dots, x_{k+1} minus ϵ ,

$$(iv) \quad \sum_{i=1}^{k+1} \omega_i^2 > 1 - \zeta.$$

We introduce some notation. The scalar product of two vectors P and Q is denoted by (P, Q) , the space spanned by the vectors P_1, P_2, \dots, P_l is denoted by $[P_1, P_2, \dots, P_l]$. For the squared length of the projection of a vector Y on some space we write the s.l.p. of Y .

Before proving the theorem it is convenient to have the following lemma's.

Lemma 1: Let P, Q, R and S be unit vectors in n -space, and let the vector Z be defined by $Z = \mu(P - \eta R) + \nu(Q - \eta S)$, where μ, ν and η are non-zero real numbers. If $(P, R) = (Q, S) = \eta$ and $(P, S) = (Q, R) = (R, S) = 0$, then the s.l.p. of Z on the space $[P, Q]$ is larger than $(Z, Z) - \eta^2(\mu^2 + \nu^2)$.

Proof: It is well known, that the s.l.p. of Z on $[P, Q]$ is given by

$$\frac{1}{1 - (P, Q)^2} \left\{ (P, \mu(P - \eta R) + \nu(Q - \eta S))^2 + (Q, \mu(P - \eta R) + \nu(Q - \eta S))^2 - 2(P, Q) \cdot (P, \mu(P - \eta R) + \nu(Q - \eta S)) \cdot (Q, \mu(P - \eta R) + \nu(Q - \eta S)) \right\}$$

$$\begin{aligned}
&= \frac{1}{1-(P,Q)^2} \{(\mu+v(P,Q)-\mu\eta^2)^2 + (v+\mu(P,Q)-v\eta^2)^2 + \\
&\quad -2(P,Q) \cdot (\mu+v(P,Q)-\mu\eta^2) \cdot (v+\mu(P,Q)-v\eta^2)\} \\
&= \mu^2 + v^2 + 2\mu v(P,Q) - 2\eta^2(\mu^2 + v^2) + \frac{\eta^4}{1-(P,Q)^2} \{\mu^2 + v^2 - 2\mu v(P,Q)\}.
\end{aligned}$$

As $(Z,Z) = \mu^2 + v^2 + 2\mu v(P,Q) - \eta^2(\mu^2 + v^2)$, and $-1 < (P,Q) < 1$, the desired result follows immediately.

Lemma 2: In the above model, disregarding the relations (3), the conditions necessary and sufficient to obtain the sequence of independent variables x_1, x_2, \dots, x_m with both forward selection and backward elimination are

$$\omega_1^2 > \omega_2^2 > \dots > \omega_{k+1}^2 \quad (6)$$

$$\omega_{k+2}^2 > \omega_{k+3}^2 > \dots > \omega_m^2 \quad (7)$$

$$\begin{aligned}
\left[\omega_i - \sum_{j=k+2}^1 \frac{\xi_{ij}}{\eta} \omega_j \right]^2 \cdot \left[1 + \sum_{j=k+2}^1 \left(\frac{\xi_{ij}}{\eta} \right)^2 \right]^{-1} > \omega_1^2 \\
(i=1, \dots, k+1; l=k+2, \dots, m) \quad (8)
\end{aligned}$$

$$\begin{aligned}
\omega_{l+1}^2 > \left[\omega_j + \sum_{i=l+1}^{k+1} \frac{\xi_{ij}}{\eta} \omega_i \right]^2 \cdot \left[1 + \sum_{i=l+1}^{k+1} \left(\frac{\xi_{ij}}{\eta} \right)^2 \right]^{-1} \\
(j=k+2, \dots, m; l=0, 1, \dots, k). \quad (9)
\end{aligned}$$

Proof: A. Backward elimination

To obtain the sequence x_1, x_2, \dots, x_m we must at each step in the process of backward elimination delete the variable with the highest index. Suppose at a certain stage we have arrived at the subset $\{x_1, x_2, \dots, x_l\}$. The condition for deleting x_l is, that the s.l.p. of Y on $[X_1, X_2, \dots, X_{l-1}]$ is larger than the s.l.p. of Y on the space spanned by the vectors $\{X_j | j=1, 2, \dots, l; j \neq l\}$ for every integer $l \leq m$.

The s.l.p. of Y on $[X_1, X_2, \dots, X_{l-1}]$ is equal to $\sum_{j=1}^{l-1} \omega_j^2$, because this space is identical with $[Y_1, Y_2, \dots, Y_{l-1}]$.

We derive the s.l.p. of Y on the second space, denoted by $v^2(i, l)$, and discuss three cases.

A1. Let $l \leq k+1$. Then $v^2(i, l) = \sum_{j=1}^l \omega_j^2 - \omega_i^2$ follows immediately.

A2. Let $k+1 < i \leq l-1$. Then clearly $v^2(i, l) = \sum_{j=1}^l \omega_j^2 - \omega_i^2$.

A3. Let $i \leq k+1 < l$. Without loss of generality we take $i=1$. Then $v^2(i, l) = (\text{s.l.p. of } Y \text{ on } [X_2, X_3, \dots, X_{k+1}]) + (\text{s.l.p. of } Y \text{ on } [\xi_{1j} Y_1 + Y_j | j=k+2, \dots, l])$. The first term of the right hand member is

equal to $\sum_{j=2}^{k+1} \omega_j^2$.

To compute the second term we define $U = \pi_1 Y_1 + \pi_{k+2} Y_{k+2} + \dots + \pi_l Y_l$ as a vector orthogonal to the second space onto which Y is projected. The orthogonality relations are $\pi_1 \xi_{1j} + \pi_j \eta = 0$ ($j=k+2, \dots, m$), from which the direction cosines of U may be computed. It follows that the distance from the endpoint of the vector $\omega_1 Y_1 + \omega_{k+2} Y_{k+2} + \dots + \omega_l Y_l$ to the space

$[\xi_{1j} Y_1 + \eta Y_j | j=k+2, \dots, l]$ is $[\omega_1 - \sum_{j=k+2}^l \frac{\xi_{1j} \omega_j}{\eta}] \cdot [1 + \sum_{j=k+2}^l (\frac{\xi_{1j}}{\eta})^2]^{-\frac{1}{2}}$.

It is now an easy matter to obtain the second term, and adding the first term we get

$$v^2(1, l) = \sum_{j=1}^l \omega_j^2 - [\omega_1 - \sum_{j=k+2}^l \frac{\xi_{1j} \omega_j}{\eta}]^2 \cdot [1 + \sum_{j=k+2}^l (\frac{\xi_{1j}}{\eta})^2]^{-1}.$$

The formulas for $v^2(2, l), \dots, v^2(k+1, l)$ are direct analogues of $v^2(1, l)$. As consequences of A1, A2 and A3 combined with the s.l.p. of Y on $[X_1, X_2, \dots, X_{l-1}]$ we find the necessary and sufficient conditions (6), (7) and (8).

B. Forward selection

Suppose at a certain step in the process of forward selection we have arrived at the subset $\{x_1, x_2, \dots, x_l\}$. To obtain the sequence x_1, x_2, \dots, x_m we must now select x_{l+1} in the regression equation. The condition is that the s.l.p. of Y on $[X_1, X_2, \dots, X_l, X_{l+1}]$ is larger than the s.l.p. of Y on $[X_1, X_2, \dots, X_l, X_j]$ for every integer $j > l+1$. The first quantity is equal to $\sum_{i=1}^{l+1} \omega_i^2$, the second, which we denote by $w^2(j, l)$, we now compute.

B1. Let $0 \leq l < k+1 < j$, where $l=0$ means that no variable has been selected yet. Then $w^2(j, l) = (\text{s.l.p. of } Y \text{ on } [X_1, X_2, \dots, X_l]) + (\text{s.l.p. of } Y \text{ on the vector } \xi_{l+1, j} Y_{l+1} + \dots + \xi_{k+1, j} Y_{k+1} + \eta Y_j) =$

$$= \sum_{i=1}^l \omega_i^2 + \left[\omega_j + \sum_{i=l+1}^{k+1} \frac{\xi_{ij}}{\eta} \omega_j \right]^2 \cdot \left[1 + \sum_{i=l+1}^{k+1} \left(\frac{\xi_{ij}}{\eta} \right)^2 \right]^{-1}.$$

B2. Let $1 \leq l+1 < j \leq k+1$. Then $w^2(j, l) = \sum_{i=1}^l \omega_i^2 + \omega_j^2$.

B3. Let $k+1 \leq l < j-1$. Then clearly $w^2(j, l) = \sum_{i=1}^l \omega_i^2 + \omega_j^2$.

We see that B2 and B3 do not lead to new conditions other than (6) or (7). As a consequence of B1 we obtain condition (9).

Proof of Theorem: Choose an η satisfying

$$0 < \eta < (1 + \max_j d_j)^{-1} \cdot \min \left\{ \frac{2}{3} \sqrt{\epsilon} \min_{i \neq j} |d_i - d_j|, \frac{\epsilon}{3} \min_j d_j \right\}. \quad (10)$$

Next choose a set of ratio's $\left\{ \frac{\omega_j}{\omega_{k+1}} \mid j=k+2, \dots, m \right\}$ satisfying

$$\frac{\omega_j}{\omega_{k+1}} < \frac{\omega_{j+1}}{\omega_{k+1}} \quad (j=k+2, \dots, m-1) \text{ and}$$

$$0 < \frac{\omega_j}{\omega_{k+1}} < \min \left\{ \frac{\sqrt{\xi}}{\sqrt{m-k-1}}, \frac{\eta d_j}{2(m-k)\sqrt{1+2/k}} \right\} \quad (j=k+2, \dots, m). \quad (11)$$

Taking $\omega_1 > 0$ and using (i) and (2) we find a set of values $\{\omega_i | i=1, \dots, m\}$, all larger than zero. Taking all $\gamma_j < 0$ and using (i) and (4) we find sets of values $\{\gamma_j | j=k+2, \dots, m\}$ and $\{\delta_j | j=k+2, \dots, m\}$. Thus property (i) is satisfied. We will show that the model with these coefficients also has the properties (ii), (iii) and (iv).
a). First we exhibit property (ii). We have to demonstrate that the inequalities (6), (7), (8) and (9) of lemma 2 are satisfied. The inequalities (6) and (7) are evident.

Consider the set of inequalities (9). Inserting (3) in (9) we find after some calculations, that (9) is equivalent to

$$\begin{aligned} \eta^2 + \gamma_j^2 \sum_{i=1+1}^{k+1} \omega_i^2 + \delta_j^2 + 2\gamma_j \delta_j \omega_{k+1} &> \eta^2 \left(\frac{\omega_j}{\omega_{1+1}} \right)^2 + 2\eta \delta_j \frac{\omega_j \omega_{k+1}}{\omega_{1+1}^2} + \\ + 2\eta \gamma_j \frac{\omega_j}{\omega_{1+1}} \sum_{i=1+1}^{k+1} \omega_i^2 + \left(\frac{\gamma_j}{\omega_{1+1}} \right)^2 \left(\sum_{i=1+1}^{k+1} \omega_i^2 \right)^2 &+ \\ + 2\gamma_j \delta_j \frac{\omega_{k+1}}{\omega_{1+1}^2} \sum_{i=1+1}^{k+1} \omega_i^2 + \delta_j^2 \left(\frac{\omega_{k+1}}{\omega_{1+1}} \right)^2 & \end{aligned}$$

for $j=k+2, \dots, m$ and $l=0, 1, \dots, k$.

We note that

$$\eta^2 > \eta^2 \left(\frac{\omega_j}{\omega_{1+1}} \right)^2 + 2\eta \delta_j \frac{\omega_j \omega_{k+1}}{\omega_{1+1}^2} \quad \text{if} \quad \frac{\omega_j}{\omega_{k+1}} < \min \left\{ \frac{1}{\sqrt{2}}, \frac{\eta}{4\delta_j} \right\}$$

and

$$\begin{aligned} \gamma_j^2 \sum_{i=1+1}^{k+1} \omega_i^2 + 2\gamma_j \delta_j \omega_{k+1} &\geq \left(\frac{\gamma_j}{\omega_{1+1}} \right)^2 \left(\sum_{i=1+1}^{k+1} \omega_i^2 \right)^2 + \\ + 2\gamma_j \delta_j \frac{\omega_{k+1}}{\omega_{1+1}^2} \sum_{i=1+1}^{k+1} \omega_i^2 &\quad \text{if} \quad \frac{|\gamma_j|}{\delta_j} < 2\omega_{k+1}; \end{aligned}$$

if both conditions are satisfied then (9) holds.

Because (11) implies $\frac{\omega_j}{\omega_{k+1}} < \frac{1}{\sqrt{m-k-1}}$ ($j=k+2, \dots, m$), and hence, using

(i), $\sum_{i=1}^{k+1} \omega_i^2 > \frac{k+1}{k+2}$, we find

$$\frac{|\gamma_j|}{\delta_j} = d_j < \left(1 + \sum_{i=1}^k c_i^2\right)^{-\frac{1}{2}} = \omega_{k+1} \left(\sum_{i=1}^{k+1} \omega_i^2\right)^{-\frac{1}{2}} < \omega_{k+1} \sqrt{\frac{k+2}{k+1}}. \quad (12)$$

Thus the second condition is satisfied.

On the other hand (4) and (11) imply

$$\delta_j^2 = (1-\eta^2) \cdot \left[\left(\frac{\gamma_j}{\delta_j}\right)^2 \sum_{i=1}^k \omega_i^2 + \left(\frac{\gamma_j}{\delta_j} \omega_{k+1} + 1\right)^2 \right]^{-1} < \frac{k+2}{k} \left(\frac{\gamma_j}{\delta_j}\right)^{-2} = \frac{1}{d_j^2} \left(1 + \frac{2}{k}\right). \quad (13)$$

Then we deduce from (11) and (13) $\frac{\omega_j}{\omega_{k+1}} < \frac{\eta d_j}{4\sqrt{1+2/k}} < \frac{\eta}{4\delta_j}$,

and the first condition is also satisfied.

We now turn to the set of inequalities (8). Inserting (3) in (8) we see, that the inequalities (8) certainly hold for $i=1, \dots, k$. For $i=k+1$ they are equivalent to

$$\eta^2 \omega_{k+1}^2 + \left(\sum_{j=k+2}^1 (\gamma_j \omega_{k+1} + \delta_j) \omega_j \right)^2 > \eta^2 \omega_1^2 + \omega_1^2 \sum_{j=k+2}^1 (\gamma_j \omega_{k+1} + \delta_j)^2 + 2\eta \omega_{k+1} \sum_{j=k+2}^1 (\gamma_j \omega_{k+1} + \delta_j) \omega_j.$$

Since from (12) $|\gamma_j| \omega_{k+1} < \delta_j \omega_{k+1} \sqrt{\frac{k+2}{k+1}} < \delta_j$, the inequalities are satisfied if

$$\eta^2 \omega_{k+1}^2 > \eta^2 \omega_1^2 + 2\eta \omega_{k+1} \sum_{j=k+2}^1 (\gamma_j \omega_{k+1} + \delta_j) \omega_j \quad (l=k+2, \dots, m).$$

Hence it is sufficient to show, that

$$\eta^2 \omega_{k+1}^2 > \eta^2 \omega_1^2 + 2\eta \omega_{k+1} \sum_{j=k+2}^1 \delta_j \omega_j.$$

From (11) we have $\omega_1^2 < \frac{1}{m-k} \omega_{k+1}^2$, and from (11) and (13)

$$\sum_{j=k+2}^1 \delta_j \frac{\omega_j}{\omega_{k+1}} < \frac{1}{2} (m-k)^{-1} \left(1 + \frac{2}{k}\right)^{-\frac{1}{2}} \eta \sum_{j=k+2}^1 d_j \delta_j < \frac{1}{2} (m-k)^{-1} (m-k-1) \eta.$$

$$\text{Thus } \eta^2 \omega_1^2 + 2\eta \omega_{k+1} \sum_{j=k+2}^1 \delta_j \omega_j < \frac{1}{m-k} \eta^2 \omega_{k+1}^2 + \frac{m-k-1}{m-k} \eta^2 \omega_{k+1}^2 = \eta^2 \omega_{k+1}^2,$$

and the inequalities (8) are satisfied.

This proves property (ii).

b). To verify property (iii) we apply lemma 1.

Consider a pair of vectors X_p, X_q , where $p, q \geq k+2$ and $p \neq q$. Inserting $P=X_p, Q=X_q, R=Y_p, S=Y_q$ and $Z=\omega_1 Y_1 + \dots + \omega_{k+1} Y_{k+1}$ in the lemma, it is

$$\text{not difficult to show that } \mu = \frac{1}{\delta_p} \left(\frac{Y_p}{\delta_p} - \frac{Y_q}{\delta_q}\right)^{-1} \text{ and } \nu = \frac{1}{\delta_q} \left(\frac{Y_q}{\delta_q} - \frac{Y_p}{\delta_p}\right)^{-1}.$$

As the conditions of the lemma are satisfied, we find that the s.l.p. of $\omega_1 Y_1 + \omega_2 Y_2 + \dots + \omega_{k+1} Y_{k+1}$ on $[\bar{X}_p, \bar{X}_q]$ is larger than

$$\sum_{i=1}^{k+1} \omega_i^2 - \eta^2 (\mu^2 + \nu^2). \text{ We investigate the last term. From (4) and (10) it}$$

is found that

$$\begin{aligned} \delta_j^2 &= (1-\eta^2) \left[\left(\frac{Y_j}{\delta_j}\right)^2 \sum_{i=1}^k \omega_i^2 + \left(\frac{Y_j}{\delta_j} \omega_{k+1} + 1\right)^2 \right]^{-1} > \frac{15}{16} \left[\left(\frac{Y_j}{\delta_j}\right)^2 + 2 \frac{|Y_j|}{\delta_j} + 1 \right]^{-1} \geq \\ &\geq \frac{15}{16} (1 + \max_j d_j)^{-2}. \end{aligned} \quad (14)$$

Using (14) and (10) we have

$$\begin{aligned} \eta^2 (\mu^2 + \nu^2) &= \eta^2 \left(\frac{1}{\delta_p^2} + \frac{1}{\delta_q^2} \right) \left(\frac{Y_p}{\delta_p} - \frac{Y_q}{\delta_q} \right)^{-2} < \frac{32}{15} \eta^2 \left[1 + \max_j d_j \right]^2 \cdot \\ &\cdot \left[\min_{i \neq j} (d_i - d_j) \right]^{-2} < \epsilon. \end{aligned}$$

Hence the s.l.p. of $\omega_1 Y_1 + \omega_2 Y_2 + \dots + \omega_{k+1} Y_{k+1}$ on $[\bar{X}_p, \bar{X}_q]$ is larger than

$\sum_{i=1}^{k+1} \omega_i^2 - \varepsilon$. A fortiori this inequality holds for the s.l.p. of Y

on $[X_p, X_q]$, and property (iii) is proved for every pair of variables from the set $\{x_l | l=k+2, \dots, m\}$.

Now consider the pair of vectors X_{k+1}, X_l , where $l > k+1$. Then the s.l.p. of Y on $[X_{k+1}, X_l]$ is equal to

$$\begin{aligned} & \omega_{k+1}^2 + \left[\sum_{i=1}^k \xi_{i1} \omega_i + \eta \omega_l \right]^2 \cdot \left[n^2 + \sum_{i=1}^k \xi_{i1}^2 \right]^{-1} = \\ & = \omega_{k+1}^2 + \left[\gamma_1 \sum_{i=1}^k \omega_i + \eta \omega_l \right]^2 \cdot \left[n^2 + \gamma_1^2 \sum_{i=1}^k \omega_i^2 \right]^{-1} = \\ & = \sum_{i=1}^{k+1} \omega_i^2 - \left[n^2 \sum_{i=1}^k \omega_i^2 + 2\eta |\gamma_1| \omega_l \sum_{i=1}^k \omega_i - \eta^2 \omega_l^2 \right] \cdot \left[n^2 + \gamma_1^2 \sum_{i=1}^k \omega_i^2 \right]^{-1} > \\ & > \sum_{i=1}^{k+1} \omega_i^2 - \frac{\eta^2}{\gamma_1^2} - \frac{\eta}{|\gamma_1|} . \end{aligned}$$

But from (i), (14) and (10) we have

$$\frac{\eta}{|\gamma_1|} = \frac{\eta}{d_1} \cdot \frac{1}{\delta_1} < \sqrt{\frac{16}{15}} \cdot \frac{\eta}{d_1} (1 + \max_j d_j) < \frac{\varepsilon}{2}, \text{ and thus } \frac{\eta^2}{\gamma_1^2} < \frac{\varepsilon^2}{4} < \frac{\varepsilon}{2} .$$

Using these inequalities we find, that the s.l.p. of Y on $[X_{k+1}, X_l]$

is larger than $\sum_{i=1}^{k+1} \omega_i^2 - \varepsilon$. This completes the proof of property (iii).

c). Property (iv) is almost trivial. From (11) we have

$$\omega_j < \sqrt{\frac{\zeta}{m-k-1}} \omega_{k+1} \quad (j=k+2, \dots, m), \text{ hence } \sum_{j=k+2}^m \omega_j^2 < \zeta, \text{ and so } \sum_{i=1}^{k+1} \omega_i^2 > 1 - \zeta.$$

Remark: The orthogonality of the vectors X_1, X_2, \dots, X_{k+1} is not an essential feature of the general example, and was only introduced to get tractable inequalities.

From the theorem we derive, that for suitable values of the coefficients in the model both forward selection and backward elimination lead to the sequence of independent variables x_1, x_2, \dots, x_m , while at the same time all k -subsets containing at least two variables from the set $\{x_j | j=k+1, \dots, m\}$ are better k -subsets than $\{x_1, x_2, \dots, x_k\}$, since their sums of squares due

to regression are larger than $\sum_{i=1}^{k+1} \omega_i^2 - \epsilon > \sum_{i=1}^k \omega_i^2$ (for a value of

ϵ smaller than $(1-\zeta) \cdot (1 + \sum_{i=1}^k c_i^2)^{-1}$). It follows that in this case

there are at least

$$\binom{m}{k} - k(m-k) - 1$$

k -subsets better than the k -subset yielded by forward selection and backward elimination. If $m=10$ and $k=3$, this means that from the totality of $\binom{10}{3}=120$ 3-subsets there are at least 98 3-subsets better than $\{x_1, x_2, x_3\}$.

In practice it is not only important how many k -subsets may exist which are better than the k -subset we choose, but also how much better they may be. We define the amount, by which a k -subset is better than another k -subset, as the difference between the sums of squares due to regression on the first and the second subset, divided by the total sum of squares. In the numerical examples discussed in section 2, the amount by which other 2-subsets were superior to $\{x_1, x_2\}$ was not very substantial. The following corollaries to our theorem show, that the situation may be much worse.

Corollary 1: For appropriate values of the coefficients in the general example, all $\binom{m}{k} - k(m-k) - 1$ k -subsets, which contain at least two variables from the set $\{x_j | j=k+1, \dots, m\}$, are better than the k -subset $\{x_1, x_2, \dots, x_k\}$ yielded by forward selection and backward elimination by an amount arbitrarily close $\frac{1}{k+1}$.

Proof: In the theorem choose $c_i^2 = 1 + \frac{2i}{k} \epsilon$ ($i=1, \dots, k$), so

$\omega_{k+1}^2 = \frac{1}{(k+1)(1+\epsilon)} \sum_{i=1}^{k+1} \omega_i^2 > \frac{1}{k+1} \cdot \frac{1-\zeta}{1+\epsilon}$. From (iii) we have, that the

amount, by which the $\binom{m}{k} - k(m-k) - 1$ k -subsets containing at least two variables from the set $\{x_l | l=k+1, \dots, m\}$ are better than $\{x_1, x_2, \dots, x_k\}$,

is larger than $\frac{\omega_{k+1}^2}{1+\theta^2} > \frac{1}{k+1} \cdot \frac{1-\zeta}{(1+\theta^2)(1+\epsilon)} - \frac{\epsilon}{1+\theta^2}$.

Because θ, ϵ and ζ may be taken arbitrarily small, the corollary follows.

Corollary 2: For appropriate values of the coefficients in the general example, all $\binom{m}{k} - (m-1) \binom{1}{k-1} - \binom{1}{k}$ k -subsets, which contain at least two variables from the set $\{x_j | j=1+1, \dots, m\}$, are better than the k -subset $\{x_1, x_2, \dots, x_k\}$ yielded by both forward selection and backward elimination by an amount arbitrarily close to $\frac{1-k+1}{1+1}$ ($k \leq 1 < m-1$).

Proof: Apply the theorem with 1 instead of k . Define

$$c_i^2 = 1 + \frac{2i}{1} \epsilon \quad (i=1, \dots, 1), \text{ so } \omega_{1+1}^2 = \frac{1}{(1+1)(1+\epsilon)} \sum_{i=1}^{1+1} \omega_i^2 > \frac{1}{1+1} \cdot \frac{1-\zeta}{1+\epsilon}.$$

Then we have from (iii), that all 2-subsets contained in the set $\{x_j | j=1+1, \dots, m\}$ are associated with sums of squares due to regression larger than

$$\sum_{i=1}^{1+1} \omega_i^2 - \epsilon > \sum_{i=1}^k \omega_i^2 + \frac{1-k+1}{1+1} \cdot \frac{1-\zeta}{1+\epsilon} - \epsilon, \text{ and the corollary follows.}$$

Corollary 3: For appropriate values of the coefficients in the general example forward selection and backward elimination, though identical, do not produce one optimal subset, except for the trivial cases of the 1-subset and $(m-1)$ -subset.

Proof: Apply the theorem with $k=m-2$ and sufficiently small ϵ .

We illustrate the corollaries with two numerical examples for the case $m=10$ and $k=3$.

From corollary 1 we know, that for suitable values of the coefficients 98 3-subsets exist better than $\{x_1, x_2, x_3\}$ by an amount close to $1/4$. In table IV a matrix of crossproduct sums is given with these properties. We only give the sums of squares, divided by the total sum of squares, due to regression on the 21 2-subsets present in the set $\{x_j | j=4, \dots, 10\}$ because the 3-subsets containing such 2-subsets can only have larger sums of squares.

In table V a numerical example is presented illustrating corollaries 2 and 3. Note that the optimal 2-subset $\{x_9, x_{10}\}$ is better than the 2-subset $\{x_1, x_2\}$ yielded by forward selection and backward elimination by an amount 0.7513. The only optimal subsets yielded by both methods are $\{x_1\}$ and $\{x_1, x_2, \dots, x_9\}$. Applying forward selection or backward elimination one would have to include 8 or 9 independent variables in the regression equation to get a good fit of the data, although the subset $\{x_9, x_{10}\}$ does the job as well.

4. Conclusion

Reviewing the results, we see that a class of examples can be constructed where forward selection and backward elimination do not lead to optimal k-subsets, even if both methods yield identical sequences of the independent variables. The k-subset they produce can be a bad one in a quantitative sense, that is, there are many better k-subsets, as well as in a qualitative sense, that is, there exists at least one k-subset that is very much better. Furthermore it is possible that both methods, though identical, do not lead to optimal k-subsets for any k except $k=1$ and $k=m-1$. In some cases it is possible to detect such anomalies by inspection of the correlation matrix: a highly intercorrelated subset of independent variables which appear in the regression equation only at a later stage, may be a sign of misbehaviour.

However, in our opinion a better and yet not too troublesome method will be hard to find, because such a method should essentially use the correlations between all variables at every stage of the process.

TABLE IV Example 4

Matrix of Crossproduct Sums

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}	y
x_1	1	0	0	0	-0.493122	-0.437261	-0.370512	-0.109394	-0.289304	-0.199584	0.501498
x_2		1	0	0	-0.491649	-0.435955	-0.369406	-0.109067	-0.288440	-0.198988	0.5
x_3			1	0	-0.491158	-0.435519	-0.369036	-0.108958	-0.288151	-0.198789	0.499500
x_4				1	0.523055	0.654814	0.767974	0.981823	0.865902	0.938490	0.498989
x_5					1	0.986372	0.947273	0.674630	0.878915	0.784771	-0.477433
x_6						1	0.986656	0.785746	0.944747	0.875133	-0.328035
x_7							1	0.875045	0.985070	0.941552	-0.171614
x_8								1	0.944666	0.986627	0.326121
x_9									1	0.985058	-0.001142
x_{10}										1	0.169433
y											1.01

Sum of squares due to Regression on the 3-subset yielded by forward selection and backward elimination as a fraction of the total sum of squares 1.01:

$$SS \{x_1, x_2, x_3\} = 0.74356.$$

Sums of Squares due to Regression on 21 different 2-subsets in fractions of the total sum of squares:

$$\begin{aligned}
 SS \{x_4, x_5\} &= 0.98974 & SS \{x_6, x_8\} &= 0.98874 & SS \{x_5, x_7\} &= 0.98520 & SS \{x_5, x_6\} &= 0.97248 \\
 SS \{x_4, x_6\} &= 0.98966 & SS \{x_4, x_{10}\} &= 0.98820 & SS \{x_7, x_{10}\} &= 0.98516 & SS \{x_9, x_{10}\} &= 0.97105 \\
 SS \{x_4, x_7\} &= 0.98950 & SS \{x_5, x_9\} &= 0.98782 & SS \{x_8, x_9\} &= 0.98507 & SS \{x_7, x_9\} &= 0.97104 \\
 SS \{x_5, x_8\} &= 0.98920 & SS \{x_7, x_8\} &= 0.98780 & SS \{x_6, x_9\} &= 0.98501 & SS \{x_8, x_{10}\} &= 0.97004 \\
 SS \{x_4, x_9\} &= 0.98917 & SS \{x_6, x_{10}\} &= 0.98778 & SS \{x_4, x_8\} &= 0.98393 & SS \{x_6, x_7\} &= 0.96993 \\
 SS \{x_5, x_{10}\} &= 0.98876 & & & & & &
 \end{aligned}$$

TABLE V Example 5

Matrix of Crossproduct Sums

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}	y
x_1	1	0	0	0	0	0	0	0	0	-0.209081	0.343511
x_2		1	0	0	0	0	0	0	0	-0.207302	0.340588
x_3			1	0	0	0	0	0	0	-0.205507	0.337639
x_4				1	0	0	0	0	0	-0.203696	0.334664
x_5					1	0	0	0	0	-0.201869	0.331662
x_6						1	0	0	0	-0.200950	0.330151
x_7							1	0	0	-0.200026	0.328634
x_8								1	0	-0.199097	0.327109
x_9									1	0.816275	0.325563
x_{10}										1	-0.278242
y											1.01

Sum of squares due to Regression on the 2-subset $\{x_9, x_{10}\}$ as a fraction of the total sum of squares 1.01:

$$SS \{x_9, x_{10}\} = 0.98298.$$

Sums of squares due to Regression on the subsets yielded by forward selection and backward elimination in fractions of the total sum of squares:

SS $\{x_1\}$	= 0.11683	SS $\{x_1, x_2, x_3, x_4, x_5\}$	= 0.56436
SS $\{x_1, x_2\}$	= 0.23168	SS $\{x_1, x_2, x_3, x_4, x_5, x_6\}$	= 0.67228
SS $\{x_1, x_2, x_3\}$	= 0.34455	SS $\{x_1, x_2, x_3, x_4, x_5, x_6, x_7\}$	= 0.77921
SS $\{x_1, x_2, x_3, x_4\}$	= 0.45545		
	SS $\{x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8\}$	= 0.88515	
	SS $\{x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9\}$	= 0.99009	
	SS $\{x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9, x_{10}\}$	= 0.99010	

References:

- [1] H.C. HAMAKER, On multiple regression analyses, *Statistica Neerlandica* 16 (1962), pp. 31-56.
- [2] P.G. MOORE, Regression as an analytical tool, *Appl.Statist.* 11 (1962), pp. 106-119.

