



DECUS

PROGRAM LIBRARY

DECUS NO.	FOCAL8-76
TITLE	SCREENING REGRESSION
AUTHOR	Dr. Robert G. Miller Glastonbury, Connecticut
COMPANY	
DATE	October 3, 1969
SOURCE LANGUAGE	FOCAL

ATTENTION

This is a USER program. Other than requiring that it conform to submittal and review standards, no quality control has been imposed upon this program by DECUS.

The DECUS Program Library is a clearing house only; it does not generate or test programs. No warranty, express or implied, is made by the contributor, Digital Equipment Computer Users Society or Digital Equipment Corporation as to the accuracy or functioning of the program or related material, and no responsibility is assumed by these parties in connection therewith.

SCREENING REGRESSION

DECUS Program Library Write-up

DECUS No. FOCAL8-76

OPERATING INSTRUCTIONS

- 1) A PDP-8 , 8/S , 8/I or 8/L with 4096 words and an ASR33 teletype.
- 2) FOCAL 1969 with no extended functions and the following suppressions or changes:

<u>IN</u>	<u>FIND</u>	<u>PUT</u>
1217	4551	7600
6002	4551	7600
2163	4551	7000
63	2676	1354
64	2666	2414
2732	6001	5336
2762	6046	7000

- 3) After reading in the FOCAL tape of "Screening Regression" start with GOTO 2.10.
- 4) Type in the correct answers to:

- P (the number of independent variables in the problem)
- N (the sample size over which the sums, sums of squares and sums of crossproducts were derived)
- SY (the sum of Y, i.e., $\sum Y$)
- SSY (the sum of squares of Y, i.e., $\sum Y^2$)
- F [the tabled F distribution value satisfying the following:]

$$F_{\alpha}(1, N) \text{ such that } = \frac{1}{20 \cdot P}$$

- 5) Place in the reader a punched tape containing the following input with P independent variables ($P=1, \dots$, no limit).

$$\Sigma X_1, \Sigma X_1^2, \Sigma X_1 X_2, \Sigma X_1 X_3, \dots, \Sigma X_1 X_p, \Sigma X_1 Y$$

$$\Sigma X_2, \Sigma X_2 X_1, \Sigma X_2^2, \Sigma X_2 X_3, \dots, \Sigma X_2 X_p, \Sigma X_2 Y$$

$$\Sigma X_3, \Sigma X_3 X_1, \Sigma X_3 X_2, \Sigma X_3^2, \Sigma X_3 X_4, \dots, \Sigma X_3 X_p, \Sigma X_3 Y$$

⋮

$$\Sigma X_p, \Sigma X_p X_1, \Sigma X_p X_2, \dots, \Sigma X_p X_{p-1}, \Sigma X_p^2, \Sigma X_p Y$$

6) Turn teletype reader to ON and wait for print of PASS 1's selected variable ID with its computed F value. If no variable is significant only E(O) will print out.

7) To run subsequent passes reenter the punched tape in the reader at the beginning. Continue doing this until the program stops selecting. This is determined when the next variable to be selected is not statistically significant or when the maximum of six variables has been achieved. The information on the seventh variable will be printed but will not be included in the regression equation.

8) The program will calculate and print the regression equation as follows:

B(O) XXX.XXXXX (Additive constant)

B(ID X⁽¹⁾) XXX.XXXXX (Coefficient for the first selected variable X⁽¹⁾)

B(ID X⁽²⁾) XXX.XXXXX (Coefficient for the second selected variable X⁽²⁾)

⋮

B(ID X^(S)) XXX.XXXXX (Coefficient for the Sth selected variable X^(S))

SCREENING REGRESSION

C-FOCAL, 1969

- $\emptyset 1.5\emptyset$ F $I=1-P, \emptyset; F J=1-P, I; S A(-J-I*P-I)=AC(-J-I*P)$
 $\emptyset 1.55$ S $B=\emptyset; F Q=1, N; D 5$
 $\emptyset 1.7\emptyset$ S $D(P)=Y; S R(P)=M; S X=X-B; F J=\emptyset, P; S A(P*J+J+P)=V(J)$
 $\emptyset 1.8\emptyset$ I $(B*A/X-W)1\emptyset.1; T \%3, P, D(P), \%8.\emptyset 2, B*A/X, !; I (6-P)1\emptyset.1; S P=P+1; G$

 $\emptyset 2.1\emptyset$ T "SCREENING REGRESSION", !
 $\emptyset 2.2\emptyset$ A "P", N, !, "N", A, !, S Y", R, !, "SS Y", X, !, "F", W, !
 $\emptyset 2.3\emptyset$ T "PASS SELVAR F", !; S P=1; S X=X-R $\uparrow 2/A; G$

 $\emptyset 3.7\emptyset$ F $J=1, P; S I=P; F L=\emptyset, J-1; D 9$
 $\emptyset 3.85$ A $R(P); S J=P; S I=P; F L=\emptyset, J-1; D 6$

 $\emptyset 4.1\emptyset$ S $Q=Q-1; G 1.7$

 $\emptyset 5.\emptyset 5$ F $J=\emptyset, P-1; D 7$
 $\emptyset 5.\emptyset 7$ S $D(P)=Q; I (N-Q)4.1; F L=\emptyset, N; A T; F J=\emptyset, P; D 8$
 $\emptyset 5.3\emptyset$ D $3; S T=R(P) \uparrow 2/A(P \uparrow 2+2*P); I (B-T)5.5; R$
 $\emptyset 5.5\emptyset$ S $B=T; S Y=Q; S M=R(P); F J=\emptyset, P; S V(J)=A(P*J+J+P)$

 $\emptyset 6.2\emptyset$ S $T=J-1-L; S T=R(T)*A(P*T+T+J)/A(P*T+2*T); S R(P)=R(P)-T$

 $\emptyset 7.1\emptyset$ I $(Q-D(J))7.2, 7.3$
 $\emptyset 7.2\emptyset$ R
 $\emptyset 7.3\emptyset$ S $Q=Q+1; F L=-1, N; A T$

 $\emptyset 8.1\emptyset$ I $(L-D(J))8.2, 8.3$
 $\emptyset 8.2\emptyset$ R
 $\emptyset 8.3\emptyset$ S $A(J+P+J*P)=T$

 $\emptyset 9.2\emptyset$ S $T=J-1-L; S T=A(I+T*P+T)*A(P*T+T+J)/A(T*P+2*T)$
 $\emptyset 9.3\emptyset$ S $A(I+P*J+J)=A(I+P*J+J)-T$

 $1\emptyset.1\emptyset$ F $I=\emptyset, P; S A(P+I*(P+1))=R(I)$
 $1\emptyset.2\emptyset$ F $J=\emptyset, P-1; F K=J-P, \emptyset; D 12$
 $1\emptyset.3\emptyset$ S $V(P-1)=A((P+1) \uparrow 2-P-2); F J=2-P, \emptyset; D 11$
 $1\emptyset.4\emptyset$ F $I=\emptyset, P-1; T \%2, "B", D(I), " ", \%6.\emptyset 5, V(I), !$
 $1\emptyset.5\emptyset$ Q

 $11.1\emptyset$ S $V(-J)=A((P+1) \uparrow 2-1-(P+1)*(P+J))$
 $11.2\emptyset$ S $T=\emptyset; F K=1, P+J-1; D 13$
 $11.3\emptyset$ S $V(-J)=V(-J)+T$

 $12.1\emptyset$ S $A(J+(P+1)*J-K)=A(J+(P+1)*J-K)/A(J+(P+1)*J)$

 $13.1\emptyset$ S $T=T-A(-J+(P+1)*(-J)+K)*V(-J+K)$

*

*SCREENING REGRESSION

P

N

S Y

SS Y

F

PASS SELVAR F

1 6 146.33

2 3 21.95

3 2 23.26

B Ø Ø.Ø8Ø16

B 6 Ø.14353

B 3 Ø.1Ø264

B 2 Ø.Ø8Ø61

*

APPENDIX A

A PROCEDURE FOR SELECTING
PREDICTORS IN MULTIPLE
REGRESSION ANALYSIS

In multiple regression a predictand, Y , is expressed as a linear function of a number of predictors, X_p ($p = 1, \dots, P$),

$$Y = a_0 + a_1X_1 + a_2X_2 + \dots + a_pX_p + \dots + a_PX_P \quad (\text{A-1})$$

where the coefficients a_p ($p = 0, \dots, P$) are determined using the method of least squares. The usual procedure for determining which of the P predictors are significant is to perform tests on the coefficients. If one or more predictors are not statistically significant, they may be eliminated from the equation by a method developed by Cochran [10]. The elimination of predictors is laborious if the original set of predictors P is large or if there are numerous non-significant predictors.

Experimenters seem to agree that in the application of multiple regression analysis to problems of prediction most of whatever predictability resides in an entire set of possible predictors can be found in a small subset of these predictors. For problems in which such is the case, it may be advantageous to attempt to select the subset of contributing predictors, leaving the redundant or noncontributing predictors out of the analysis. A method given by Bryan [7] attempts to perform just such a function. A description of the method will now be presented using the notation of Chapter I.

In order to select the first predictor $X^{(1)}$ the following quantities must be computed:

$$\begin{aligned} SST(Y) &= \sum_{i=1}^N (Y_i - \bar{Y})^2 \\ SST(X_p) &= \sum_{i=1}^N (X_{pi} - \bar{X}_p)^2 \\ SPT(YX_p) &= \sum_{i=1}^N (Y_i - \bar{Y})(X_{pi} - \bar{X}_p) \end{aligned} \quad (\text{A-2})$$

$(p = 1, \dots, P)$

where the number of observations in the dependent sample is N . Selection of the first predictor $X^{(1)}$, from the set of P possible predictors, is determined

using the criterion¹

$$\frac{SSF(X^{(1)})}{SSR(X^{(1)})} \geq \frac{SSF(X_p)}{SSR(X_p)} \quad (p = 1, \dots, P) \quad (\text{A-3})$$

where $SSF(X_p)$, the fitted sum of squares for predictor X_p on Y , is

$$SSF(X_p) = SPT(YX_p) \cdot [SST(X_p)]^{-1} \cdot SPT(YX_p). \quad (\text{A-4})$$

and $SSR(X_p)$, the residual sum of squares after fitting Y with predictor X_p , is got by subtracting $SSF(X_p)$ from $SST(Y)$, the total sum of squares of Y , or

$$SSR(X_p) = SST(Y) - SSF(X_p). \quad (\text{A-5})$$

The significance of $X^{(1)}$ may not be tested by the usual F ratio since the criterion is a function of the test statistic. An approximate test of significance for $X^{(1)}$ has been suggested [23] which uses a critical F value whose probability level is a function of the number of possible predictors P . This critical F value, F_{α^*} , is defined as

$$F_{\alpha^*} \equiv F_{\alpha} \quad (\text{A-6})$$

where α^* is the size desired in the selection test, and α is the corresponding size of the tabled F distribution. For a fixed size α^* , the value of α is taken to be approximately equal to α^*/P . The critical value of F_{α^*} for testing the significance of the S th selected predictor $X^{(S)}$ becomes

$$F_{\alpha^*} = F_{\alpha^*/(P-S+1)}. \quad (\text{A-7})$$

Therefore, the predictor $X^{(1)}$ is considered significant if

$$(N-2) \frac{SSF(X^{(1)})}{SSR(X^{(1)})} > F_{(\alpha^*/P)}(1, N-2). \quad (\text{A-8})$$

Should $X^{(1)}$ be significant it then becomes the first of r predictors to be selected from the original set of P possible predictors. If $X^{(1)}$ is not significant, no predictors are selected.

¹ In the rare event that two or more of the P predictors satisfy the inequality in (A-3), the procedure is to select, arbitrarily, $X^{(1)}$ to be the first of these in the order 1, \dots , P . This procedure is to be followed for selection of subsequent predictors when more than one predictor satisfies the selection criterion. One further point, quotients are actually not necessary in (A-3) for determining $X^{(1)}$. It would be sufficient computationally to use the criterion $SSF(X^{(1)}) \geq SSF(X_p)$. The quotient form is given to show the parallel between regression and discriminant analysis selection.

The choice of the second selected predictor $X^{(2)}$ requires the quantities

$$SPT(X^{(1)}X_p) = \sum_{i=1}^N (X_{i(1)} - \bar{X}_{(1)})(X_{pi} - \bar{X}_p) \quad (p = 1, \dots, P; X_p \neq X^{(1)}). \quad (A-9)$$

The criterion for choosing $X^{(2)}$ is

$$\frac{SSF(X^{(2)}|X^{(1)})}{SSR(X^{(2)}|X^{(1)})} \geq \frac{SSF(X_p|X^{(1)})}{SSR(X_p|X^{(1)})} \quad (p = 1, \dots, P; X_p \neq X^{(1)}) \quad (A-10)$$

where

$$SSF(X_p|X^{(1)}) = \left[\begin{matrix} SPT(YX^{(1)}) \\ SPT(YX_p) \end{matrix} \right]' \left[\begin{matrix} SST(X^{(1)}) & SPT(X^{(1)}X_p) \\ SPT(X^{(1)}X_p) & SST(X_p) \end{matrix} \right]^{-1} \left[\begin{matrix} SPT(YX^{(1)}) \\ SPT(YX_p) \end{matrix} \right] - SSF(X^{(1)}) \quad (p = 1, \dots, P; X_p \neq X^{(1)}) \quad (A-11)$$

and

$$SSR(X_p|X^{(1)}) = SST(Y) - SSF(X^{(1)}) - SSF(X_p|X^{(1)}) \quad (p = 1, \dots, P). \quad (A-12)$$

Predictor $X^{(2)}$ is judged significant if, from (A-7),

$$(N-3) \frac{SSF(X^{(2)}|X^{(1)})}{SSR(X^{(2)}|X^{(1)})} > F_{\alpha*/(P-1)}(1, N-3). \quad (A-13)$$

Selection of predictor $X^{(S)}$ is made using the following general form:

$$\frac{SSF(X^{(S)}|X^{(1)} \dots X^{(S-1)})}{SSR(X^{(S)}|X^{(1)} \dots X^{(S-1)})} \geq \frac{SSF(X_p|X^{(1)} \dots X^{(S-1)})}{SSR(X_p|X^{(1)} \dots X^{(S-1)})} \quad (p = 1, \dots, P; X_p \neq X^{(1)} \dots X^{(S-1)}) \quad (A-14)$$

where

$$SSF(X_p|X^{(1)} \dots X^{(S-1)}) = \left[\begin{matrix} SPT(YX^{(1)}) \\ \vdots \\ SPT(YX^{(S-1)}) \\ SPT(YX_p) \end{matrix} \right]' \left[\begin{matrix} SST(X^{(1)}) & \dots & SPT(X^{(1)}X^{(S-1)}) & SPT(X^{(1)}X_p) \\ \vdots & & \vdots & \vdots \\ SPT(X^{(1)}X^{(S-1)}) & \dots & SST(X^{(S-1)}) & SPT(X^{(S-1)}X_p) \\ SPT(X^{(1)}X_p) & \dots & SPT(X^{(S-1)}X_p) & SST(X_p) \end{matrix} \right]^{-1} \left[\begin{matrix} SPT(YX^{(1)}) \\ \vdots \\ SPT(YX^{(S-1)}) \\ SPT(YX_p) \end{matrix} \right] - SSF(X^{(1)}) - \dots - SSF(X^{(S-1)}|X^{(1)} \dots X^{(S-2)}) \quad (p = 1, \dots, P; X_p \neq X^{(1)} \dots X^{(S-1)}) \quad (A-15)$$

and

$$SSR(X_p|X^{(1)} \dots X^{(S-1)}) = SST(Y) - SSF(X^{(1)}) - SSF(X^{(2)}|X^{(1)}) - \dots - SSF(X_p|X^{(1)} \dots X^{(S-1)}) \quad (p = 1, \dots, P; X_p \neq X^{(1)} \dots X^{(S-1)}). \quad (A-16)$$

Predictor $X^{(S)}$ is said to be significant if

$$[N - (S+1)] \frac{SSF(X^{(S)}|X^{(1)} \dots X^{(S-1)})}{SSR(X^{(S)}|X^{(1)} \dots X^{(S-1)})} > F_{\alpha*/(P-S+1)}(1, N - (S+1)). \quad (A-17)$$

For a particular regression problem the number of selected predictors r is determined such that the predictor $X^{(r+1)}$ fails to show significance.

This procedure cannot be said to select necessarily the best set of r predictors out of the original set containing P predictors. However, it has been shown that it can select a highly reliable set of predictors when applied to particular problems in meteorology [24; 27; 35].

This method may leave a set of unselected predictors which together contain significant predictive

information. Because none of these alone gives a significant contribution they all remain unselected. This can be remedied in part by testing significance on combinations of selected variables. In practice, however, it has been found that F^* tends to work well as a significance test only when variables are considered singly. This may be a consequence of the fact that a bias is introduced in the regression coefficients as a

result of selection. A critical evaluation of the amount of bias created by selection has not as yet been carried out.

The method contains the labor-saving feature of not requiring the crossproducts between unselected predictors. When all P variables must be processed

in ordinary multiple regression $\frac{P(P+1)}{2}$ sums of squares or crossproducts are required. In the selection procedure only $P \cdot r - \frac{r(r-1)}{2}$ are necessary. If P is large this saving is sizable

