

DOCUMENT RESUME

ED 171 746

TE 008 824

AUTHOR Pohlmann, John T.
 TITLE Controlling the Type I Error Rate in Stepwise Regression Analysis.
 PUE DATE Apr 79
 NOTE 21p.; Paper presented at the Annual Meeting of the American Educational Research Association (63rd, San Francisco, California, April 8-12, 1979); appendices marginally legible

EDRS PRICE MF01/PC01 Plus Postage.
 DESCRIPTORS Computer Programs; *Error Patterns; *Mathematical Models; *Multiple Regression Analysis; Predictor Variables; Research Design; Statistical Analysis; *Statistical Studies; *Tests of Significance

IDENTIFIERS *Stepwise Regression; *Type I Errors

ABSTRACT

Three procedures used to control Type I error rate in stepwise regression analysis are forward selection, backward elimination, and true stepwise. In the forward selection method, a model of the dependent variable is formed by choosing the single best predictor; then the second predictor which makes the strongest contribution to the prediction of the dependent variable is chosen, controlling for the effects of the first variable. The process continues so that the variable chosen increases the prediction potential, until remaining variables fail to make any contribution. Backward elimination begins with a model containing all predictors; and, at each step, a variable is eliminated if its removal results in the smallest reduction of effectiveness. True stepwise procedure is a variant of forward selection. To test these procedures, a Monte Carlo computer program, written in FORTRAN IV, was prepared. The results support two conclusions: (1) the probability of erroneously forming a regression model increases as a function of the number of predictors; and (2) as the inter-predictor correlation increases, the probability of making errors decreases. Therefore, the number of predictors and the inter-predictor correlation should be considered when attempting to solve an error rate problem. (MH)

 * Reproductions supplied by EDRS are the best that can be made *
 * from the original document. *

ED171746

U.S. DEPARTMENT OF HEALTH,
EDUCATION & WELFARE
NATIONAL INSTITUTE OF
EDUCATION

THIS DOCUMENT HAS BEEN REPRODUCED EXACTLY AS RECEIVED FROM THE PERSON OR ORGANIZATION ORIGINATING IT. POINTS OF VIEW OR OPINIONS STATED DO NOT NECESSARILY REPRESENT OFFICIAL NATIONAL INSTITUTE OF EDUCATION POSITION OR POLICY

Controlling the Type I Error Rate
in Stepwise Regression Analysis

John T. Pohlmann
Southern Illinois University, Carbondale

"PERMISSION TO REPRODUCE THIS
MATERIAL HAS BEEN GRANTED BY

John T. Pohlmann

TO THE EDUCATIONAL RESOURCES
INFORMATION CENTER (ERIC) AND
USERS OF THE ERIC SYSTEM."

AERA

April, 1979

San Francisco

Session 24.19

TM008 824

Controlling the Type I Error Rate in Stepwise Regression Analysis

Stepwise regression has become a widely used technique for selecting a subset of potential predictors for some dependent variable. Three procedures have been used under the rubric of stepwise regression analysis: forward selection, backward elimination, and true stepwise (Draper and Smith, 1966).

The forward selection procedure forms a model of the dependent variable by first selecting the best single predictor, then the second predictor is chosen to maximize the marginal contribution to the prediction of Y , controlling for the effect of the first predictor. The process continues until at each step, the variable selected for inclusion in the model increases the prediction of Y more than any other predictor. The selection process stops when the remaining variables fail to contribute significantly to the prediction of Y . The backward elimination procedure begins with a model containing all potential predictors, and then at each step a variable is eliminated if its removal from the model results in the smallest reduction in the model's effectiveness. The elimination process continues until the removal of any variable results in a significant reduction in the model's R^2 . The true stepwise procedure is a variation of the forward selection technique. It differs from the forward selection procedure in that at each step, a variable

that has been previously included in the model may be deleted if a partial F-test shows that variable to be an insignificant predictor.

In most of the computer statistical packages that have stepwise regression procedures, the criterion used for variable selection is an F-test formed as follows:

$$F = \frac{R_F^2 - R_R^2}{(1 - R_F^2) / (N - p - 1)} \quad (1)$$

where: R_F^2 = the coefficient of determination for the model containing all predictors included at previous steps, plus the variable under test.

R_R^2 = the coefficient of determination for the model containing all predictors except the variable under test.

N = the number of observations.

p = the number of predictors used in the model that produced R_F^2 .

As with any statistical test, two kinds of inferential errors can be made. A type I error could occur if a variable was selected, using the F ratio criterion, when that variable's population regression weight was zero. A type II error occurs when a variable is not selected, using the F-test criterion, when that variable has a non-zero population weight.

Most users of stepwise regression adopt one of the traditional

significance levels (.05 or .01) when evaluating the F-test in (1). This significance level will determine the type I error rate for each test. However, another perspective can be taken when considering the type I error rate, the problem-wide error rate.

The problem-wide error rate is the probability of selecting any variable when all variables have population regression weights of zero. In other words, the problem-wide error rate is the probability of forming a sample regression model, when none should be formed. The rest of this paper addresses this error rate, and a procedure will be presented that allows researchers to control its value.

The problem-wide error rate is comparable to the family-wide error rate commonly encountered in the context of post hoc tests conducted after a significant effect has been found in an ANOVA. For example, the probability of making one or more type I errors in a family of orthogonal tests is:

$$\alpha_F = 1 - \prod_{i=1}^k (1 - \alpha_i) \quad (2)$$

where α_F = the family-wide error rate.
 k = the number of orthogonal tests.
 α_i = the significance level on test i .

When the α_i 's are all equal to α_T ,

$$\alpha_F = 1 - (1 - \alpha_T)^k. \quad (3)$$

If a researcher wished to control α_F by reducing α_T , (3) could be solved for α_T :

$$\alpha_T = 1 - \sqrt[k]{1 - \alpha_F}. \quad (4)$$

Alternately, the researcher could conservatively approximate α_T using the Bonferroni inequality,

$$\alpha_T \approx \alpha_F/k \quad (5)$$

When the members of the family of tests are not orthogonal, formulas (4) and (5) yield conservative values of α_T . That is, the use of α_T from (4) or (5) will result in an α_F less than the desired value. The solution for α_T is considerably more complex when the tests are not orthogonal. The solution for a critical F that will maintain α_F at a desired value should be done using the correlated F distribution (Pope and Webster, 1972). Unfortunately the integration of the correlated F distribution is an extremely tedious process, and only limited tables critical values derived from it are available. Consequently, an approximate solution was sought using Monte Carlo methods.

METHOD

A Monte Carlo program, written in FORTRAN IV, was prepared by the author for this project. The program incorporated subroutines supplied in the International Mathematical and Statistical Library (1975). The IMSL subroutines were selected because of their proven accuracy and efficiency. A copy of the program is supplied in the Appendix of this paper.

The program generated sample data matrices (cases by variables) sampled with a given population dispersion matrix. Subroutine GGNRM was used for this purpose. Various population correlation matrices were supplied to GGNRM and a sample data matrix of standard normal deviates was produced. All population correlations between the predictors and the criterion variable were set equal to zero. The inter-predictor correlations were all set equal to a common value, and for the various replications examined in this study, the inter-predictor correlations were .2, .3, .5, .7, and .9. In addition, the numbers of predictors used were 2, 3, 4, 5, 10, and 20. For every combination of the number of predictors and the average inter-predictor correlation (35 in all), a thousand sample data sets were generated.

Each data set thus generated was then subjected to a stepwise regression analysis using IMSL subroutine RLSTEP. Subroutine RLSTEP uses a true stepwise procedure. Variable selection is governed by a significance testing process. When, at any step, no F-test is significant, the selection process ceases.

For the purposes of this study, an error occurred when a model, other than the null model, was formed by subroutine RLSTEP. The proportion of analyses resulting in a model was treated as an empirical estimate of the probability of erroneously forming a model using stepwise regression analysis.

RESULTS

Table 1 shows the results obtained when a variable selection significance level of .05 is used. The table entries in Table 1 are the proportion of 1000 stepwise regression analyses that produced a sample

model when none should have been produced. For example, when a researcher has ten potential predictors that have correlations with each other equal to .50, the probability of erroneously forming a model is approximately .308.

Since the values in Table 1 are empirical estimates of the actual probabilities of making an error, there is some sampling error. The magnitude of ~~the~~ sampling error can be conservatively estimated by using the standard error of a proportion when $p = .5$. Since 1000 replications were used to derive each table entry, the standard error of a sample proportion will be less than or equal to .016. Consequently, a conservative 68% confidence interval for the true probability of making an error will be: tabulated value \pm .016.

The figures in Table 1 suggest two conclusions: (1) The probability of erroneously forming a regression model increases dramatically as a function of the number of predictors, and (2), as the inter-predictor correlation increases, the probability of making an error decreases. Consequently, any solution to the error rate problem must take into consideration the number of predictors and the inter-predictor correlation.

After Table 1 was prepared, an attempt was made to develop an algorithm that could be used to select a significance level for variable selection that would control the problem-wide error rate.

The rationale for the algorithm presented here was based on the formula that gives the family-wide error rate in k independent tests. Formula (3) is reproduced here for this purpose:

$$\alpha_F = 1 - (1 - \alpha_T)^k, \quad (6)$$

All terms are defined in (3). If α_T and α_F are known, k can be solved for as follows:

$$k = \frac{\ln(1-\alpha_F)}{\ln(1-\alpha_T)} \quad (7)$$

Formula (7) was applied to each entry in Table 1, and the resulting k values are given in Table 2. In producing Table 2, α_T was .05 and α_F was taken as the corresponding value in Table 1. The k values in Table 2 were then plotted as a function of various measures of the inter-predictor correlation. Figure 1 shows one of these plots for the 10 predictor variable case. The k values were observed to be an inverse linear function of ρ_{XX}^2 , the inter-predictor correlation. The following function was considered to be a reasonable approximation:

$$k = p - (p-1)\rho_{XX}^2 \quad (8)$$

where

p = the number of predictors

ρ_{XX}^2 = the inter-predictor correlation.

This function seemed suitable since for the extreme cases of ρ_{XX}^2 , 0 and 1.0, (8) produced k values of p and 1 respectively. When ρ_{XX}^2 is equal to 0, the problem-wide error rate should equal the α_F value given by (6). Under this condition ($\rho_{XX}^2 = 0$) the error rate is directly analogous to the family-wide error rate for a family of orthogonal tests. When ρ_{XX}^2 is equal to 1, every predictor is linearly dependent on the other predictors, hence there is in fact only one predictor. Formula (8) yields a k value of 1, when ρ_{XX}^2 equals 1. In addition, inspection of plots, such

as Figure 1, suggested that (8) was also accurate for estimating k for values of ρ_{xx}^2 between 0 and 1.

Unfortunately a researcher using stepwise regression never knows ρ_{xx}^2 , so it must be estimated. A less biased estimate of the squared correlation coefficient can be obtained using the shrinkage formula (McNemar, 1969):

$$\hat{\rho}^2 = 1 - (1-r^2) \frac{N-1}{N-2} \quad . \quad (9)$$

The estimate of ρ_{xx}^2 used for this study was obtained as follows:

Let R_{pp} = the inter-predictor correlation matrix.

Define each element in \hat{R}_{pp} as

$$\hat{r}_{ij}^2 = 1 - (1-r_{ij}^2) \frac{N-1}{N-2} \quad , \quad (10)$$

where r_{ij}^2 = the square of the ij th element of R_{pp} , and

N = the number of observations.

$$\text{Let } \bar{r}^2 = \frac{\sum_{i=1}^{p-1} \sum_{j=i+1}^p \hat{r}_{ij}^2}{\frac{1}{2}(p^2-p)} \quad , \quad (11)$$

which is the mean of the off diagonal elements of \hat{R}_{pp} .

The sample estimate of ρ_{xx}^2 is then substituted into (8) to obtain

$$k = p - (p-1)\bar{r}^2 \quad . \quad (12)$$

After k has been obtained via (12), α_T is obtained.

$$\alpha_T = 1 - \sqrt[k]{1 - \alpha_F} \quad , \quad (13)$$

where α_F is the desired problem-wide error rate.

A concise worked example is given in the Appendix of this paper.

The validity of the proposed algorithm was then tested by modifying the Monte Carlo program, used to produce Table 1, to use (13) to select an α_T . The results of this validation study are presented in Table 3. As can be noted in Table 3, the probability of erroneously forming a model, using (13) to determine α_T , approaches the desired value of .05. There is a slight tendency for this procedure to produce conservative values of α_T . The average value of α_F in Table 3 is .045, and the conservative nature of the procedure is most apparent for problems with large numbers of predictors and high inter-predictor correlations.

DISCUSSION

The type I error rate in stepwise regression analysis deserves serious consideration by researchers. The literature is replete with "significant" findings that fail the ultimate test of replication. One possible explanation for this state of affairs might lie in the increasing problem-wide error rate that can occur in stepwise regression analysis.

If a researcher considers the problem wide error rate important, he or she should take some corrective action. Three possibilities exist, depending on the kind of analysis contemplated. They are: (1) Prior to the stepwise analysis conduct an omnibus test of the model containing all potential predictors, (2) use the backward elimination procedure and use an α_T obtained by substituting the number of predictors for k in (13), or (3) use the algorithm for obtaining α_T presented here, if a forward

selection or true stepwise procedure is used.

The Omnibus Test

The analysis begins by forming a full model containing all predictors. The R^2 for this model is tested for significance at the α_F level. The F is obtained as follows:

$$F = \frac{R^2/p}{(1-R^2)/(N-p-1)} \quad , \quad (14)$$

where R^2 = the coefficient of determination for the model
 containing all potential predictors,
 p = the number of predictors,
 N = the number of cases.

This F ratio yields a simultaneous test of significance for all weights in a model. Proceed with the analysis only if a significant F using (14) is obtained.

The Backward Elimination Procedure

The backward elimination procedure is comparable to testing a family of orthogonal hypotheses. At each step, the variance accounted for in the dependent variable that is tested for each predictor is independent of all the sources of variation. Consequently, the use of

$$\alpha_T = 1 - \sqrt[p]{1-\alpha_F} \quad , \quad (15)$$

will maintain α_F at its desired value.

Finally, the algorithm developed in this paper is recommended if a forward selection or true stepwise procedure is used. Since the value of α_T obtained using (13) will be greater than that obtained using (15),

when some covariance among the predictors is present, the use of (13) will produce a more powerful analysis.

References

- Draper, N. R., & Smith, H. Applied regression analysis. New York: John Wiley & Sons, 1966.
- International mathematical and statistical library (5th Ed.). Houston: International Mathematical and Statistical Libraries, Inc., 1975.
- McNemar, Q. Psychological statistics. New York: John Wiley & Sons, 1969.
- Pope, P. T., & Webster, J. T. The use of an F-statistic in stepwise regression procedures. Technometrics, 1972, 14(2), 327-340.

Table 1
 Monte Carlo Estimates of the Probability of
 Erroneously Forming a Sample Model Using
 Stepwise Regression Analysis with
 a Variable Selection Significance Level of .05

Inter-Predictor Correlation	Number of Predictors						
	2	3	4	5	7	10	20
.0	.102	.130	.184	.216	.304	.410	.653
.3	.101	.130	.178	.213	.275	.367	.552
.5	.097	.128	.171	.196	.235	.308	.417
.7	.085	.125	.140	.153	.185	.225	.314
.9	.073	.094	.101	.111	.122	.126	.169

Table 2
k Values Derived Using Formula (7)
on the Values from Table 1

Inter-Predictor Correlation	Number of Predictors						
	2	3	4	5	7	10	20
.0	2.10	2.72	3.96	4.74	7.06	10.29	20.63
.3	2.08	2.72	3.82	4.67	6.27	8.92	15.70
.5	1.99	2.67	3.66	4.25	5.22	7.18	10.52
.7	1.73	2.60	2.94	3.24	3.98	4.97	7.35
.9	1.48	1.92	2.08	2.28	2.54	2.63	3.61

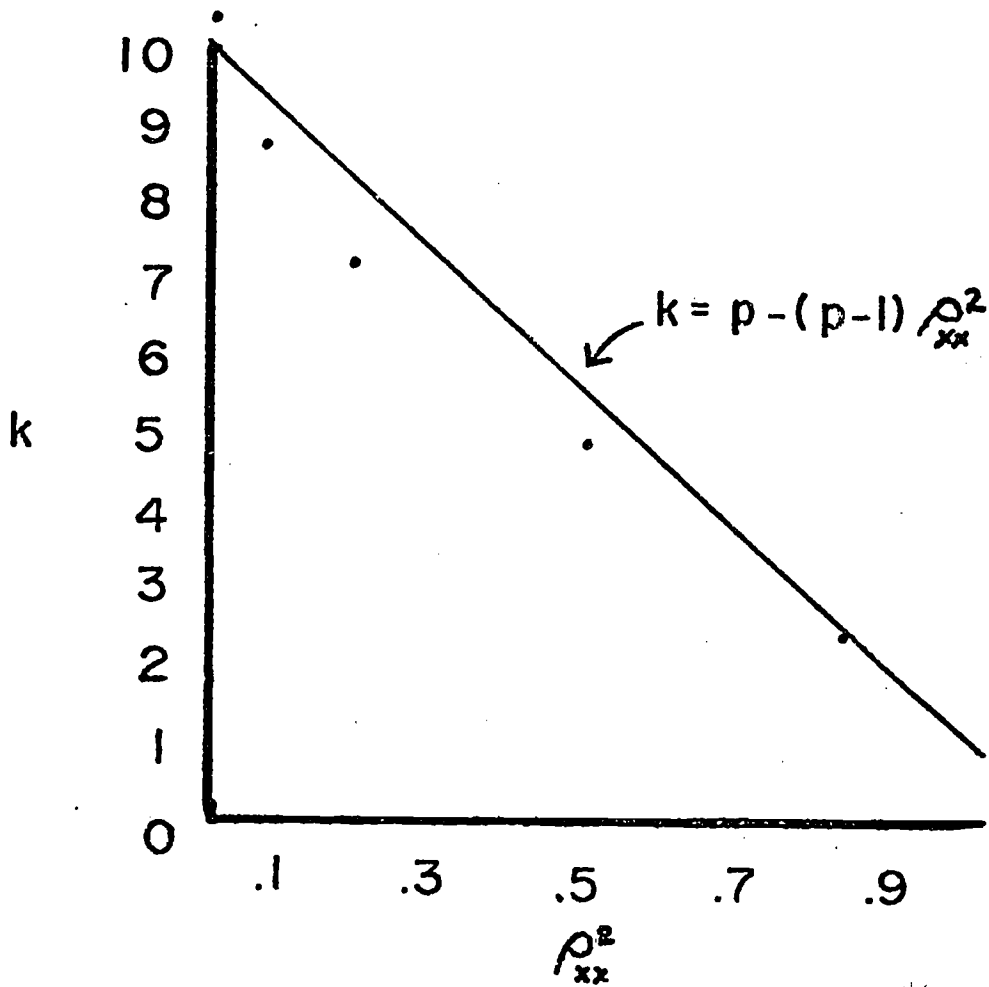


Figure 1. Plot of k as a function of ρ_{xx}^2 for 10 predictors.

Table 3
 Monte Carlo Estimates of the Probability
 of Erroneously Forming a Sample Model Using Stepwise
 Regression Analysis with a Variable Selection Significance
 Level Obtained Using Formula 13. The Desired α_F was .05

Inter-Predictor Correlation	Number of Predictors						
	2	3	4	5	7	10	20
.0	.052	.044	.058	.044	.048	.045	.055
.3	.050	.045	.044	.055	.046	.047	.038
.5	.060	.044	.041	.063	.041	.044	.042
.7	.059	.050	.041	.046	.037	.031	.032
.9	.045	.054	.056	.050	.033	.027	.011

APPENDIX I

Copy of the Computer Program Used
in the Study

```

1000 CALL RND(500)
1010 CALL RND(1000)/RND(1000)/RND(1000),I1(11)
1020 I1(11)=
1030 WRITE(6,15)
1040 CALL IAI('INPUT SEED')
1050 READ(5,*) ISEED
1060 WRITE(6,16)
1070 FORMAT(1X,'INPUT SEEDS, VALUES AGAIN: I1(11)')
1080 READ(5,*) I1(11)=999999,K,ALPHA,ALPOUT
1090 WRITE(6,20)
1100 CALL('INPUT LINE NUMBER OF RND: I1(11) IS I1(11)')
1110 I1(11)=I1(11)/2
1120 READ(5,*) I1(11)=999999,I1(11),I1(11)
1130 WRITE(6,20)
1140 CALL IAI('INPUT SEEDS, VALUES AGAIN: I1(11)')
1150 READ(5,*) I1(11),I1(11)
1160 I1(11)=0
1170 I1(11)=I1(11)*I1(11),NSAMP
1180 WRITE(6,302) I1(11)
1190 I1(11)=I1(11)*I1(11),NSAMP
1200 I1(11)=I1(11)*I1(11),NSAMP
1210 I1(11)=I1(11)*I1(11),NSAMP
1220 I1(11)=I1(11)*I1(11),NSAMP
1230 I1(11)=I1(11)*I1(11),NSAMP
1240 I1(11)=I1(11)*I1(11),NSAMP
1250 I1(11)=I1(11)*I1(11),NSAMP
1260 I1(11)=I1(11)*I1(11),NSAMP
1270 I1(11)=I1(11)*I1(11),NSAMP
1280 I1(11)=I1(11)*I1(11),NSAMP
1290 I1(11)=I1(11)*I1(11),NSAMP
1300 I1(11)=I1(11)*I1(11),NSAMP
1310 I1(11)=I1(11)*I1(11),NSAMP
1320 I1(11)=I1(11)*I1(11),NSAMP
1330 I1(11)=I1(11)*I1(11),NSAMP
1340 I1(11)=I1(11)*I1(11),NSAMP
1350 I1(11)=I1(11)*I1(11),NSAMP
1360 I1(11)=I1(11)*I1(11),NSAMP
1370 I1(11)=I1(11)*I1(11),NSAMP
1380 I1(11)=I1(11)*I1(11),NSAMP
1390 I1(11)=I1(11)*I1(11),NSAMP
1400 I1(11)=I1(11)*I1(11),NSAMP
1410 I1(11)=I1(11)*I1(11),NSAMP
1420 I1(11)=I1(11)*I1(11),NSAMP
1430 I1(11)=I1(11)*I1(11),NSAMP
1440 I1(11)=I1(11)*I1(11),NSAMP
1450 I1(11)=I1(11)*I1(11),NSAMP
1460 I1(11)=I1(11)*I1(11),NSAMP
1470 I1(11)=I1(11)*I1(11),NSAMP
1480 I1(11)=I1(11)*I1(11),NSAMP
1490 I1(11)=I1(11)*I1(11),NSAMP
1500 I1(11)=I1(11)*I1(11),NSAMP
1510 I1(11)=I1(11)*I1(11),NSAMP
1520 I1(11)=I1(11)*I1(11),NSAMP
1530 I1(11)=I1(11)*I1(11),NSAMP
1540 I1(11)=I1(11)*I1(11),NSAMP
1550 I1(11)=I1(11)*I1(11),NSAMP
1560 I1(11)=I1(11)*I1(11),NSAMP
1570 I1(11)=I1(11)*I1(11),NSAMP
1580 I1(11)=I1(11)*I1(11),NSAMP
1590 I1(11)=I1(11)*I1(11),NSAMP
1600 I1(11)=I1(11)*I1(11),NSAMP
1610 I1(11)=I1(11)*I1(11),NSAMP
1620 I1(11)=I1(11)*I1(11),NSAMP
1630 I1(11)=I1(11)*I1(11),NSAMP
1640 I1(11)=I1(11)*I1(11),NSAMP
1650 I1(11)=I1(11)*I1(11),NSAMP
1660 I1(11)=I1(11)*I1(11),NSAMP
1670 I1(11)=I1(11)*I1(11),NSAMP
1680 I1(11)=I1(11)*I1(11),NSAMP
1690 I1(11)=I1(11)*I1(11),NSAMP
1700 I1(11)=I1(11)*I1(11),NSAMP
1710 I1(11)=I1(11)*I1(11),NSAMP
1720 I1(11)=I1(11)*I1(11),NSAMP
1730 I1(11)=I1(11)*I1(11),NSAMP
1740 I1(11)=I1(11)*I1(11),NSAMP
1750 I1(11)=I1(11)*I1(11),NSAMP
1760 I1(11)=I1(11)*I1(11),NSAMP
1770 I1(11)=I1(11)*I1(11),NSAMP
1780 I1(11)=I1(11)*I1(11),NSAMP
1790 I1(11)=I1(11)*I1(11),NSAMP
1800 I1(11)=I1(11)*I1(11),NSAMP
1810 I1(11)=I1(11)*I1(11),NSAMP
1820 I1(11)=I1(11)*I1(11),NSAMP
1830 I1(11)=I1(11)*I1(11),NSAMP
1840 I1(11)=I1(11)*I1(11),NSAMP
1850 I1(11)=I1(11)*I1(11),NSAMP
1860 I1(11)=I1(11)*I1(11),NSAMP
1870 I1(11)=I1(11)*I1(11),NSAMP
1880 I1(11)=I1(11)*I1(11),NSAMP
1890 I1(11)=I1(11)*I1(11),NSAMP
1900 I1(11)=I1(11)*I1(11),NSAMP
1910 I1(11)=I1(11)*I1(11),NSAMP
1920 I1(11)=I1(11)*I1(11),NSAMP
1930 I1(11)=I1(11)*I1(11),NSAMP
1940 I1(11)=I1(11)*I1(11),NSAMP
1950 I1(11)=I1(11)*I1(11),NSAMP
1960 I1(11)=I1(11)*I1(11),NSAMP
1970 I1(11)=I1(11)*I1(11),NSAMP
1980 I1(11)=I1(11)*I1(11),NSAMP
1990 I1(11)=I1(11)*I1(11),NSAMP
2000 I1(11)=I1(11)*I1(11),NSAMP

```



APPENDIX II

A Worked Example of the Algorithm for Obtaining a Significance Level for Variable Selection Using Stepwise Regression

R_{pp}

$N = 20$

Desired Model Error Rate = .05

$$\begin{bmatrix} 1.0 & .3 & .5 \\ & 1.0 & .2 \\ \text{sym} & & 1.0 \end{bmatrix}$$

$$\hat{r}_{ij}^2 = 1 - (1 - r_{ij}^2) \frac{N - 1}{N - 2}$$

$$\hat{r}_{12}^2 = .0394$$

$$\hat{r}_{13}^2 = .2083$$

$$\hat{r}_{23}^2 = -.0133$$

$$\bar{r}^2 = \frac{\sum_{i=1}^{p-1} \sum_{j=i+1}^p \hat{r}_{ij}^2}{\frac{1}{2}(p^2 - p)} = .0781$$

$$k = p - (p - 1) \bar{r}^2 = 2.8438$$

$$\alpha_T = 1 - \sqrt[k]{1 - \alpha_F} = .0179$$