Stepwise Regression: nagdmc_stepwise_reg

Purpose

 ${\bf nagdmc_stepwise_reg}$ computes a forward stepwise regression on p variables using Clarke's sweep algorithm.

Declaration

Parameters

1:	rec1 – long On entry: the index in the data of the first data record used in the analysis. Constraint: rec1 ≥ 0 .	Input
2:	nvar – long On entry: the number of variables in the data. Constraint: $nvar > 1$.	Input
3:	nrec – long On entry: the number of consecutive records, beginning at rec1 , used in the analysis. Constraint: $nrec > 1$.	Input
4:	$dblk - long$ $On entry:$ the total number of records in the data block. $Constraint:$ $dblk \geq rec1 + nrec.$	Input
5:	data[dblk * nvar] - double On entry: the data values for the <i>j</i> th variable (for $j = 0, 1,, nvar-1$) are stored in $data[i*nvar]$ for $i = 0, 1,, dblk - 1$.	$Input \\ \mathbf{r+}j],$
6:	nxvar – long Input On entry: the number of independent variables. If nxvar = 0 then all variables in the data, excluding yvar and iwts , are treated as independent variables. Constraint: $0 \leq nxvar < nvar$.	
7:	xvar [nxvar] - long Input On entry: the indices indicating the position in data in which values of the independent variables are stored. If nxvar = 0 then xvar must be 0, and the indices of independent variables are given by $j = 0, 1,, $ nvar - 1; $j \neq$ yvar and $i \neq$ iwts . Constraints: if nxvar > 0, $0 \leq$ xvar [i] < nvar , for $i = 0, 1,,$ nxvar - 1; otherwise xvar must be 0.	
8:	yvar – long On entry: the index in data in which values of the dependent variable are stored. Constraints: $0 \leq$ yvar $<$ nvar; if nxvar > 0 , yvar \neq xvar $[i]$, for $i = 0, 1,,$ nxvar -1 .	Input
9:	iwts - long On entry: if iwts = -1, no weights are used; otherwise iwts is the index in data in whice weights are stored. Constraints: $-1 \leq iwts < nvar$; $iwts \neq yvar$; and if $nxvar > 0$, $iwts \neq xvar[i]$ i = 0, 1,, nxvar - 1.	Input h the , for

Input

Input

Input

xtype[p-1] - int10:

> On entry: set xtype[j] = 2 to force the *j*th independent variable to be selected for a model; otherwise $\mathbf{xtype}[j]$ must be set equal to zero, for $j = 0, 1, \dots, p-2$, where p-1 is the number of independent variables in the analysis.

On exit: the value of xtype[j] indicates the status of the *j*th independent variable in the data for a model:

 $\mathbf{xtype}[j] = 0$; omitted $\mathbf{xtype}[j] = 1$; selected $\mathbf{xtype}[j] = 2$; forced selection

Constraint: $xtype[j] \in \{0, 2\}$, for j = 0, 1, ..., p - 2.

11: $\mathbf{Fin} - \texttt{double}$

On entry: the value of the F-test which the coefficient of a variable must exceed for that variable to be included in a model. Typical values for Fin lie in the interval [2, 4].

Constraint: Fin > 0.0.

Suggested value: $\mathbf{Fin} = 2$.

12: Fout - double

On entry: the variable in a model corresponding to the coefficient with the lowest F-test value is removed from the model if its coefficient's F-test is less than Fout. The value of Fout is usually set equal to the value of Fin; a value less than Fin is occasionally prefered.

Constraint: $0.0 < Fout \le Fin$.

Suggested value: Fout = 2.

13:eps - double

On entry: the tolerance for detecting collinearities between variables when adding or removing a variable from a model. Variables deemed to be collinear are excluded from the final model.

Constraint: eps > 0.0.

Suggested value: $eps = 1.0 \times 10^{-6}$.

14: $\mathbf{b}[p] - \mathtt{double}$

On exit: the parameter estimates. $\mathbf{b}[0]$ is the mean parameter. $\mathbf{b}[i]$ is the coefficient of the *i*th variable included in the model, for $i = 1, 2, \ldots, p-1$. If **nxvar** > 0 then the order the independent variables are added to the model is defined by **xvar**, otherwise the order is defined by indices in the data.

15:	$\mathbf{se}[p] - \mathtt{double}$	Output
	On exit: the standard errors of the parameters in \mathbf{b} .	
16:	${f R2}-{f double}$ *	Output
	On exit: the R^2 -value for the fitted model.	
17:	${f rms}-{f double}$ *	Output
	On exit: the residual mean square for the fitted model.	
18:	df - long *	Output
	On exit: the degrees of freedom for the for the residual mean square.	
19:	$\mathbf{model}[(3*p*(p+1))/2 + \mathbf{nvar} + 14] - \mathtt{double}$	Output
	On exit: if not 0, information on the fitted model for use in the functions described in 'S	ee Also'.
20:	$\mathbf{printlevel} - \mathtt{int}$	Input
	On entry: controls the information printed by the function. If printlevel is set equal to) zero, no

information is printed; otherwise **printlevel** must be set equal to one, and information is printed to a location determined by the value of file.

Constraint: **printlevel** $\in \{0, 1\}$.

Input

Output

21: file[] - char

On entry: if **printlevel** = 1, the value of **file** determines the destination of the printed summary information; otherwise **file** is not referenced. If **file** is equal to 0, information is printed to the screen of the terminal used to call **nagdmc_stepwise_reg**; otherwise information is printed to the ASCII file named **file**.

Constraint: if printlevel = 1 and file $\neq 0$, nagdmc_stepwise_reg returns an error if a file named file cannot be opened for writing.

22: info - int *

On exit: info gives information on the success of the function call:

-1:

- 0: the function successfully completed its task..one or more of the variance terms were zero when computing the correlation matrix
- $i;\,i=1,2,\ldots,4,6,7,\ldots,13,20,21$: the specification of the $i{\rm th}$ formal parameter was incorrect.
- 55: a negative value for a weight was found.
- 99: the function failed to allocate enough memory.
- 100: an internal error occurred during the execution of the function.

Notation

the number of data records, n .
determines the number of independent variables, $p-1$.
the independent variables that take the values in X .
the dependent variable that takes the values in y .
if $iwts \ge 0$, $iwts$ is the index in the data that defines the weights, W.
the minimum value, α_1 , of partial <i>F</i> -tests for entry into a model.
the minimum value, α_2 , of partial <i>F</i> -tests for removal from a model.

Description

The general multiple linear regression model is defined by

 $y = X\beta + \varepsilon,$

where y is a vector of n observations on the dependent variable,

- X is an n by p matrix of the independent variables of column rank k,
- β is a vector of length p of unknown parameters,
- and ε is a vector of length *n* of unknown, Normally distributed, random errors such that var $\varepsilon = V\sigma^2$, where *V* is a known diagonal matrix of size *n*.

In most linear regression models the first term is taken as a mean term or an intercept, i.e., $X_{i,1} = 1$, for i = 1, 2, ..., n; this is assumed in the NAG DMC.

The least squares estimates $\hat{\beta}$ of the parameters β minimise $(y - X\beta)^T (y - X\beta)$ while the weighted least squares estimates minimise $(y - X\beta)^T W(y - X\beta)$.

The goal of stepwise regression is to determine a model by selecting statistically significant variables from the p-1 available variables (the intercept is always included in a model). The forward stepwise method used to select variables incorporates an entry step and a removal step.

In the entry step the entry statistic is computed for each variable eligible for entry in the model. If no variable has a partial F-test value which exceeds the specified critical value, α_1 , then the process is terminated; otherwise the variable with the largest value on the entry statistic is entered into the model.

In the removal step the removal statistic is computed for each variable eligible to be removed from the model. If no variable has a partial F-test value which is less than a critical value, α_2 , then the process is terminated; otherwise the variable with the smallest value on the removal statistic is removed from the model.

Input

Output

The forward stepwise procedure begins by performing forward entry for all independent variables in a model. At any subsequent step where two or more variables have been selected for entry into the model, the entry procedure followed by the removal procedure is performed. The forward stepwise selection procedure halts when neither procedure can be performed.

References and Further Reading

Clarke M R B (1981) Algorithm AS 178: the Gauss-Jordan sweep operator with detection of collinearity *Applied Statistics* **31** 166-169.

Dempster A P (1969) Elements of Continuous Multivariate Analysis Addison-Wesley.

Draper N R and Smith H (1985) Applied Regression Analysis (2nd Edition) Wiley.

See Also

nagdmc_basic_reg	simplified version of nagdmc_linear_reg using a restricted set of parameters.
nagdmc_extr	computes fitted values, residuals and leverages for a regression.
nagdmc_linear_reg	linear model with Normal errors.
nagdmc_predict_reg	computes predictions given a fitted regression model.
$stepwise_reg_ex.c$	the example calling program.