

NAG Library Routine Document

D02SAF

Note: before using this routine, please read the Users' Note for your implementation to check the interpretation of ***bold italicised*** terms and other implementation-dependent details.

1 Purpose

D02SAF solves a two-point boundary value problem for a system of first-order ordinary differential equations with boundary conditions, combined with additional algebraic equations. It uses initial value techniques and a modified Newton iteration in a shooting and matching method.

2 Specification

```
SUBROUTINE D02SAF (P, M, N, N1, PE, PF, E, DP, NPOINT, SWP, LDSWP,      &
                  ICOUNT, RANGE, BC, FCN, EQN, CONSTR, YMAX, MONIT,  &
                  PRSOL, W, LDW, SDW, IFAIL)

INTEGER            M, N, N1, NPOINT, LDSWP, ICOUNT, LDW, SDW, IFAIL
REAL (KIND=nag_wp) P(M), PE(M), PF(M), E(N), DP(M), SWP(LDSWP,6), YMAX,      &
                  W(LDW,SDW)
LOGICAL           CONSTR
EXTERNAL          RANGE, BC, FCN, EQN, CONSTR, MONIT, PRSOL
```

3 Description

D02SAF solves a two-point boundary value problem for a system of n first-order ordinary differential equations with separated boundary conditions by determining certain unknown arguments p_1, p_2, \dots, p_m . (There may also be additional algebraic equations to be solved in the determination of the arguments and, if so, these equations are defined by EQN.) The arguments may be, but need not be, boundary values; they may include eigenvalues, arguments in the coefficients of the differential equations, coefficients in series expansions or asymptotic expansions for boundary values, the length of the range of definition of the system of differential equations, etc.

It is assumed that we have a system of n differential equations of the form

$$y' = f(x, y, p), \quad (1)$$

where $p = (p_1, p_2, \dots, p_m)^T$ is the vector of arguments, and that the derivative f is evaluated by FCN. Also, n_1 of the equations are assumed to depend on p . For $n_1 < n$ the $n - n_1$ equations of the system are not involved in the matching process. These are the driving equations; they should be independent of p and of the solution of the other n_1 equations. In numbering the equations in FCN and BC the driving equations must be put **first** (as they naturally occur in most applications). The range of definition $[a, b]$ of the differential equations is defined by RANGE and may depend on the arguments p_1, p_2, \dots, p_m (that is, on p). RANGE must define the points $x_1, x_2, \dots, x_{\text{NPOINT}}$, $\text{NPOINT} \geq 2$, which must satisfy

$$a = x_1 < x_2 < \dots < x_{\text{NPOINT}} = b \quad (2)$$

(or a similar relationship with all the inequalities reversed).

If $\text{NPOINT} > 2$ the points $x_1, x_2, \dots, x_{\text{NPOINT}}$ can be used to break up the range of definition. Integration is restarted at each of these points. This means that the differential equations (1) can be defined differently in each sub-interval $[x_i, x_{i+1}]$, for $i = 1, 2, \dots, \text{NPOINT} - 1$. Also, since initial and maximum integration step sizes can be supplied on each sub-interval (via the array SWP), you can indicate parts of the range $[a, b]$ where the solution $y(x)$ may be difficult to obtain accurately and can take appropriate action.

The boundary conditions may also depend on the arguments and are applied at $a = x_1$ and $b = x_{\text{NPOINT}}$. They are defined (in BC) in the form

$$y(a) = g_1(p), > y(b) = g_2(p). \quad (3)$$

The boundary value problem is solved by determining the unknown arguments p by a shooting and matching technique. The differential equations are always integrated from a to b with initial values $y(a) = g_1(p)$. The solution vector thus obtained at $x = b$ is subtracted from the vector $g_2(p)$ to give the n_1 residuals $r_1(p)$, ignoring the first $n - n_1$, driving equations. Because the direction of integration is always from a to b , it is unnecessary, in BC, to supply values for the first $n - n_1$ boundary values at b , that is the first $n - n_1$ components of g_2 in (3). For $n_1 < m$ then $r_1(p)$. Together with the $m - n_1$ equations defined by EQN,

$$r_2(p) = 0, \quad (4)$$

these give a vector of residuals r , which at the solution, p , must satisfy

$$\begin{pmatrix} r_1(p) \\ r_2(p) \end{pmatrix} = 0. \quad (5)$$

These equations are solved by a pseudo-Newton iteration which uses a modified singular value decomposition of $J = \frac{\partial r}{\partial p}$ when solving the linear equations which arise. The Jacobian J used in Newton's method is obtained by numerical differentiation. The arguments at each Newton iteration are accepted only if the norm $\|D^{-1}\tilde{J}^+r\|_2$ is much reduced from its previous value. Here \tilde{J}^+ is the pseudo-inverse, calculated from the singular value decomposition, of a modified version of the Jacobian J (J^+ is actually the inverse of the Jacobian in well-conditioned cases). D is a diagonal matrix with

$$d_{ii} = \max(|p_i|, \text{PF}(i)) \quad (6)$$

where PF is an array of floor values.

See Deuffhard (1974) for further details of the variants of Newton's method used, Gay (1976) for the modification of the singular value decomposition and Gladwell (1979) for an overview of the method used.

Two facilities are provided to prevent the pseudo-Newton iteration running into difficulty. First, you are permitted to specify constraints on the values of the arguments p via a CONSTR. These constraints are only used to prevent the Newton iteration using values for p which would violate them; that is, they are not used to determine the values of p . Secondly, you are permitted to specify a maximum value y_{\max} for $\|y(x)\|_{\infty}$ at all points in the range $[a, b]$. It is intended that this facility be used to prevent machine 'overflow' in the integrations of equation (1) due to poor choices of the arguments p which might arise during the Newton iteration. When using this facility, it is presumed that you have an estimate of the likely size of $\|y(x)\|_{\infty}$ at all points $x \in [a, b]$. y_{\max} should then be chosen rather larger (say by a factor of 10) than this estimate.

You are strongly advised to supply a MONIT (or to call the 'default' routine D02HBX, see MONIT) to monitor the progress of the pseudo-Newton iteration. You can output the solution of the problem $y(x)$ by supplying a suitable PRSOL (an example is given in Section 10 of a routine designed to output the solution at equally spaced points).

D02SAF is designed to try all possible options before admitting failure and returning to you. Provided the routine can start the Newton iteration from the initial point p it will exhaust all the options available to it (though you can override this by specifying a maximum number of iterations to be taken). The fact that all its options have been exhausted is the only error exit from the iteration. Other error exits are possible, however, whilst setting up the Newton iteration and when computing the final solution.

If you require more background information about the solution of boundary value problems by shooting methods you are recommended to read the appropriate chapters of Hall and Watt (1976), and for a detailed description of D02SAF Gladwell (1979) is recommended.

4 References

Deuffhard P (1974) A modified Newton method for the solution of ill-conditioned systems of nonlinear equations with application to multiple shooting *Numer. Math.* **22** 289–315

Gay D (1976) On modifying singular values to solve possibly singular systems of nonlinear equations *Working Paper 125* Computer Research Centre, National Bureau for Economics and Management Science, Cambridge, MA

Gladwell I (1979) The development of the boundary value codes in the ordinary differential equations chapter of the NAG Library *Codes for Boundary Value Problems in Ordinary Differential Equations. Lecture Notes in Computer Science* (eds B Childs, M Scott, J W Daniel, E Denman and P Nelson) **76** Springer–Verlag

Hall G and Watt J M (ed.) (1976) *Modern Numerical Methods for Ordinary Differential Equations* Clarendon Press, Oxford

5 Arguments

- 1: P(M) – REAL (KIND=nag_wp) array *Input/Output*
On entry: P(*i*) must be set to an estimate of the *i*th argument, p_i , for $i = 1, 2, \dots, m$.
On exit: the corrected value for the *i*th argument, unless an error has occurred, when it contains the last calculated value of the argument.
- 2: M – INTEGER *Input*
On entry: m , the number of arguments.
Constraint: $M > 0$.
- 3: N – INTEGER *Input*
On entry: n , the total number of differential equations.
Constraint: $N > 0$.
- 4: N1 – INTEGER *Input*
On entry: n_1 , the number of differential equations active in the matching process. The active equations must be placed last in the numbering in FCN and BC. The **first** $N - N1$ equations are used as the driving equations.
Constraint: $N1 \leq N$, $N1 \leq M$ and $N1 > 0$.
- 5: PE(M) – REAL (KIND=nag_wp) array *Input*
On entry: PE(*i*), for $i = 1, 2, \dots, m$, must be set to a positive value for use in the convergence test in the *i*th argument p_i . See the description of PF for further details.
Constraint: PE(*i*) > 0.0 , for $i = 1, 2, \dots, m$.
- 6: PF(M) – REAL (KIND=nag_wp) array *Input/Output*
On entry: PF(*i*), for $i = 1, 2, \dots, m$, should be set to a ‘floor’ value in the convergence test on the *i*th argument p_i . If PF(*i*) ≤ 0.0 on entry then it is set to the small positive value $\sqrt{\epsilon}$ (where ϵ may in most cases be considered to be **machine precision**); otherwise it is used unchanged.
The Newton iteration is presumed to have converged if a full Newton step is taken (ISTATE = 1 in the specification of MONIT), the singular values of the Jacobian are not being significantly perturbed (also see MONIT) and if the Newton correction C_i satisfies

$$|C_i| \leq \text{PE}(i) \times \max(|p_i|, \text{PF}(i)), \quad i = 1, 2, \dots, m,$$

where p_i is the current value of the i th argument. The values $PF(i)$ are also used in determining the Newton iterates as discussed in Section 3, see equation (6).

On exit: the values actually used.

- 7: $E(N)$ – REAL (KIND=nag_wp) array *Input*

On entry: values for use in controlling the local error in the integration of the differential equations. If err_i is an estimate of the local error in y_i , for $i = 1, 2, \dots, n$, then

$$|err_i| \leq E(i) \times \max\{\sqrt{\epsilon}, |y_i|\},$$

where ϵ may in most cases be considered to be **machine precision**.

Suggested value: $E(i) = 10^{-5}$.

Constraint: $E(i) > 0.0$, for $i = 1, 2, \dots, N$.

- 8: $DP(M)$ – REAL (KIND=nag_wp) array *Input/Output*

On entry: a value to be used in perturbing the argument p_i in the numerical differentiation to estimate the Jacobian used in Newton's method. If $DP(i) = 0.0$ on entry, an estimate is made internally by setting

$$DP(i) = \sqrt{\epsilon} \times \max(PF(i), |p_i|), \quad (7)$$

where p_i is the initial value of the argument supplied by you and ϵ may in most cases be considered to be **machine precision**. The estimate of the Jacobian, J , is made using forward differences, that is for each i , for $i = 1, 2, \dots, m$, p_i is perturbed to $p_i + DP(i)$ and the i th column of J is estimated as

$$(r(p_i + DP(i)) - r(p_i))/DP(i)$$

where the other components of p are unchanged (see (3) for the notation used). If this fails to produce a Jacobian with significant columns, backward differences are tried by perturbing p_i to $p_i - DP(i)$ and if this also fails then central differences are used with p_i perturbed to $p_i + 10.0 \times DP(i)$. If this also fails then the calculation of the Jacobian is abandoned. If the Jacobian has not previously been calculated then an error exit is taken. If an earlier estimate of the Jacobian is available then the current argument set, p_i , for $i = 1, 2, \dots, m$, is abandoned in favour of the last argument set from which useful progress was made and the singular values of the Jacobian used at the point are modified before proceeding with the Newton iteration. You are recommended to use the default value $DP(i) = 0.0$ unless you have prior knowledge of a better choice. If any of the perturbations described are likely to lead to an unfortunate set of argument values then you should use the LOGICAL FUNCTION CONSTR to prevent such perturbations (all changes of arguments are checked by a call to CONSTR).

On exit: the values actually used.

- 9: $NPOINT$ – INTEGER *Input*

On entry: 2 plus the number of break-points in the range of definition of the system of differential equations (1).

Constraint: $NPOINT \geq 2$.

- 10: $SWP(LDSWP, 6)$ – REAL (KIND=nag_wp) array *Input/Output*

On entry: $SWP(i, 1)$ must contain an estimate for an initial step size for integration across the i th sub-interval $[X(i), X(i+1)]$, for $i = 1, 2, \dots, NPOINT - 1$, (see RANGE). $SWP(i, 1)$ should have the same sign as $X(i+1) - X(i)$ if it is nonzero. If $SWP(i, 1) = 0.0$, on entry, a default value for the initial step size is calculated internally. This is the recommended mode of entry.

$SWP(i, 3)$ must contain a lower bound for the modulus of the step size on the i th sub-interval $[X(i), X(i+1)]$, for $i = 1, 2, \dots, NPOINT - 1$. If $SWP(i, 3) = 0.0$ on entry, a very small default value is used. By setting $SWP(i, 3) > 0.0$ but smaller than the expected step sizes (assuming you

have some insight into the likely step sizes) expensive integrations with arguments p far from the solution can be avoided.

$\text{SWP}(i, 2)$ must contain an upper bound on the modulus of the step size to be used in the integration on $[X(i), X(i+1)]$, for $i = 1, 2, \dots, \text{NPOINT} - 1$. If $\text{SWP}(i, 2) = 0.0$ on entry no bound is assumed. This is the recommended mode of entry unless the solution is expected to have important features which might be ‘missed’ in the integration if the step size were permitted to be chosen freely.

On exit: $\text{SWP}(i, 1)$ contains the initial step size used on the last integration on $[X(i), X(i+1)]$, for $i = 1, 2, \dots, \text{NPOINT} - 1$, (excluding integrations during the calculation of the Jacobian).

$\text{SWP}(i, 2)$, for $i = 1, 2, \dots, \text{NPOINT} - 1$, is usually unchanged. If the maximum step size $\text{SWP}(i, 2)$ is so small or the length of the range $[X(i), X(i+1)]$ is so short that on the last integration the step size was not controlled in the main by the size of the error tolerances $E(i)$ but by these other factors, then $\text{SWP}(\text{NPOINT}, 2)$ is set to the floating-point value of i if the problem last occurred in $[X(i), X(i+1)]$. Any results obtained when this value is returned as nonzero should be viewed with caution.

$\text{SWP}(i, 3)$, for $i = 1, 2, \dots, \text{NPOINT} - 1$, are unchanged.

If an error exit with $\text{IFAIL} = 4, 5$ or 6 (see Section 6) occurs on the integration made from $X(i)$ to $X(i+1)$ the floating-point value of i is returned in $\text{SWP}(\text{NPOINT}, 1)$. The actual point $x \in [X(i), X(i+1)]$ where the error occurred is returned in $\text{SWP}(1, 5)$ (see also the specification of W). The floating-point value of NPOINT is returned in $\text{SWP}(\text{NPOINT}, 1)$ if the error exit is caused by a call to BC .

If an error exit occurs when estimating the Jacobian matrix ($\text{IFAIL} = 7, 8, 9, 10, 11$ or 12 , see Section 6) and if argument p_i was the cause of the failure then on exit $\text{SWP}(\text{NPOINT}, 1)$ contains the floating-point value of i .

$\text{SWP}(i, 4)$ contains the point $X(i)$, for $i = 1, 2, \dots, \text{NPOINT}$, used at the solution p or at the final values of p if an error occurred.

SWP is also partly used as workspace.

11: LDSWP – INTEGER *Input*

On entry: the first dimension of the array SWP as declared in the (sub)program from which D02SAF is called.

Constraint: $\text{LDSWP} \geq \text{NPOINT}$.

12: ICOUNT – INTEGER *Input*

On entry: an upper bound on the number of Newton iterations. If $\text{ICOUNT} = 0$ on entry, no check on the number of iterations is made (this is the recommended mode of entry).

Constraint: $\text{ICOUNT} \geq 0$.

13: RANGE – SUBROUTINE, supplied by the user. *External Procedure*

RANGE must specify the break-points x_i , for $i = 1, 2, \dots, \text{NPOINT}$, which may depend on the arguments p_j , for $j = 1, 2, \dots, m$.

The specification of RANGE is:

```
SUBROUTINE RANGE (X, NPOINT, P, M)
```

```
  INTEGER          NPOINT, M
```

```
  REAL (KIND=nag_wp) X(NPOINT), P(M)
```

```
1:   X(NPOINT) – REAL (KIND=nag_wp) array
```

Output

On exit: the i th break-point, for $i = 1, 2, \dots, \text{NPOINT}$. The sequence $(X(i))$ must be strictly monotonic, that is either

	$a = X(1) < X(2) < \dots < X(NPOINT) = b$	
	or	
	$a = X(1) > X(2) > \dots > X(NPOINT) = b.$	
2:	NPOINT – INTEGER	<i>Input</i>
	<i>On entry:</i> 2 plus the number of break-points in (a, b) .	
3:	P(M) – REAL (KIND=nag_wp) array	<i>Input</i>
	<i>On entry:</i> the current estimate of the i th argument, for $i = 1, 2, \dots, m$.	
4:	M – INTEGER	<i>Input</i>
	<i>On entry:</i> m , the number of arguments.	

RANGE must either be a module subprogram USED by, or declared as EXTERNAL in, the (sub) program from which D02SAF is called. Arguments denoted as *Input* must **not** be changed by this procedure.

- 14: BC – SUBROUTINE, supplied by the user. *External Procedure*

BC must place in G1 and G2 the boundary conditions at a and b respectively.

The specification of BC is:		
SUBROUTINE BC (G1, G2, P, M, N)		
INTEGER M, N		
REAL (KIND=nag_wp) G1(N), G2(N), P(M)		
1:	G1(N) – REAL (KIND=nag_wp) array	<i>Output</i>
	<i>On exit:</i> the value of $y_i(a)$, for $i = 1, 2, \dots, n$, (where this may be a known value or a function of the parameters p_j , for $j = 1, 2, \dots, m$).	
2:	G2(N) – REAL (KIND=nag_wp) array	<i>Output</i>
	<i>On exit:</i> the value of $y_i(b)$, for $i = 1, 2, \dots, n$, (where these may be known values or functions of the parameters p_j , for $j = 1, 2, \dots, m$). If $n > n_1$, so that there are some driving equations, then the first $n - n_1$ values of G2 need not be set since they are never used.	
3:	P(M) – REAL (KIND=nag_wp) array	<i>Input</i>
	<i>On entry:</i> an estimate of the i th argument, p_i , for $i = 1, 2, \dots, m$.	
4:	M – INTEGER	<i>Input</i>
	<i>On entry:</i> m , the number of arguments.	
5:	N – INTEGER	<i>Input</i>
	<i>On entry:</i> n , the number of differential equations.	

BC must either be a module subprogram USED by, or declared as EXTERNAL in, the (sub) program from which D02SAF is called. Arguments denoted as *Input* must **not** be changed by this procedure.

- 15: FCN – SUBROUTINE, supplied by the user. *External Procedure*

FCN must evaluate the functions f_i (i.e., the derivatives y'_i), for $i = 1, 2, \dots, n$.

The specification of FCN is:

```
SUBROUTINE FCN (X, Y, F, N, P, M, I)
INTEGER          N, M, I
REAL (KIND=nag_wp) X, Y(N), F(N), P(M)
```

- | | | |
|----|---|---------------|
| 1: | X – REAL (KIND=nag_wp) | <i>Input</i> |
| | <i>On entry:</i> x , the value of the argument. | |
| 2: | Y(N) – REAL (KIND=nag_wp) array | <i>Input</i> |
| | <i>On entry:</i> y_i , for $i = 1, 2, \dots, n$, the value of the argument. | |
| 3: | F(N) – REAL (KIND=nag_wp) array | <i>Output</i> |
| | <i>On exit:</i> the derivative of y_i , for $i = 1, 2, \dots, n$, evaluated at x . $F(i)$ may depend upon the parameters p_j , for $j = 1, 2, \dots, m$. If there are any driving equations (see Section 3) then these must be numbered first in the ordering of the components of F. | |
| 4: | N – INTEGER | <i>Input</i> |
| | <i>On entry:</i> n , the number of equations. | |
| 5: | P(M) – REAL (KIND=nag_wp) array | <i>Input</i> |
| | <i>On entry:</i> the current estimate of the i th argument p_i , for $i = 1, 2, \dots, m$. | |
| 6: | M – INTEGER | <i>Input</i> |
| | <i>On entry:</i> m , the number of arguments. | |
| 7: | I – INTEGER | <i>Input</i> |
| | <i>On entry:</i> specifies the sub-interval $[x_i, x_{i+1}]$ on which the derivatives are to be evaluated. | |

FCN must either be a module subprogram USED by, or declared as EXTERNAL in, the (sub) program from which D02SAF is called. Arguments denoted as *Input* must **not** be changed by this procedure.

- 16: EQN – SUBROUTINE, supplied by the NAG Library or the user. *External Procedure*

EQN is used to describe the additional algebraic equations to be solved in the determination of the parameters, p_i , for $i = 1, 2, \dots, m$. If there are no additional algebraic equations (i.e., $m = n_1$) then EQN is never called and the dummy routine D02HBZ should be used as the actual argument.

The specification of EQN is:

```
SUBROUTINE EQN (E, Q, P, M)
INTEGER          Q, M
REAL (KIND=nag_wp) E(Q), P(M)
```

- | | | |
|----|---|---------------|
| 1: | E(Q) – REAL (KIND=nag_wp) array | <i>Output</i> |
| | <i>On exit:</i> the vector of residuals, $r_2(p)$, that is the amount by which the current estimates of the arguments fail to satisfy the algebraic equations. | |
| 2: | Q – INTEGER | <i>Input</i> |
| | <i>On entry:</i> the number of algebraic equations, $m - n_1$. | |

3:	P(M) – REAL (KIND=nag_wp) array	<i>Input</i>
	<i>On entry:</i> the current estimate of the i th argument p_i , for $i = 1, 2, \dots, m$.	
4:	M – INTEGER	<i>Input</i>
	<i>On entry:</i> m , the number of arguments.	

EQN must either be a module subprogram USED by, or declared as EXTERNAL in, the (sub) program from which D02SAF is called. Arguments denoted as *Input* must **not** be changed by this procedure.

- 17: CONSTR – LOGICAL FUNCTION, supplied by the user. *External Procedure*

CONSTR is used to prevent the pseudo-Newton iteration running into difficulty. CONSTR should return the value .TRUE. if the constraints are satisfied by the parameters p_1, p_2, \dots, p_m . Otherwise CONSTR should return the value .FALSE.. Usually the dummy function D02HBY, which returns the value .TRUE. at all times, will suffice and in the first instance this is recommended as the actual argument.

The specification of CONSTR is:		
FUNCTION CONSTR (P, M)		
LOGICAL CONSTR		
INTEGER M		
REAL (KIND=nag_wp) P(M)		
1:	P(M) – REAL (KIND=nag_wp) array	<i>Input</i>
	<i>On entry:</i> an estimate of the i th argument, p_i , for $i = 1, 2, \dots, m$.	
2:	M – INTEGER	<i>Input</i>
	<i>On entry:</i> m , the number of arguments.	

CONSTR must either be a module subprogram USED by, or declared as EXTERNAL in, the (sub)program from which D02SAF is called. Arguments denoted as *Input* must **not** be changed by this procedure.

- 18: YMAX – REAL (KIND=nag_wp) *Input*

On entry: a non-negative value which is used as a bound on all values $\|y(x)\|_\infty$ where $y(x)$ is the solution at any point x between X(1) and X(NPOINT) for the current arguments p_1, p_2, \dots, p_m . If this bound is exceeded the integration is terminated and the current arguments are rejected. Such a rejection will result in an error exit if it prevents the initial residual or Jacobian, or the final solution, being calculated. If YMAX = 0 on entry, no bound on the solution y is used; that is the integrations proceed without any checking on the size of $\|y\|_\infty$.

- 19: MONIT – SUBROUTINE, supplied by the NAG Library or the user. *External Procedure*

MONIT enables you to monitor the values of various quantities during the calculation. It is called by D02SAF after every calculation of the norm $\|D^{-1}\tilde{J}^+r\|_2$ which determines the strategy of the Newton method, every time there is an internal error exit leading to a change of strategy, and before an error exit when calculating the initial Jacobian. Usually the routine D02HBX will be adequate and you are advised to use this as the actual argument for MONIT in the first instance. (In this case a call to X04ABF must be made before the call of D02SAF.) If no monitoring is required, the dummy routine D02SAS may be used.

The specification of MONIT is:		
SUBROUTINE MONIT (ISTATE, IFLAG, IFAIL1, P, M, F, PNORM, PNORM1, &		
EPS, D)		

INTEGER ISTATE, IFLAG, IFAIL1, M REAL (KIND=nag_wp) P(M), F(M), PNORM, PNORM1, EPS, D(M)		
1:	ISTATE – INTEGER <i>On entry:</i> the state of the Newton iteration. ISTATE = 0 The calculation of the residual, Jacobian and $\ D^{-1}\tilde{J}^+r\ _2$ are taking place. ISTATE = 1 to 5 During the Newton iteration a factor of $2^{(-ISTATE+1)}$ of the Newton step is being used to try to reduce the norm. ISTATE = 6 The current Newton step has been rejected and the Jacobian is being re-calculated. ISTATE = -6 to -1 An internal error exit has caused the rejection of the current set of argument values, p . -ISTATE is the value which ISTATE would have taken if the error had not occurred. ISTATE = -7 An internal error exit has occurred when calculating the initial Jacobian.	<i>Input</i>
2:	IFLAG – INTEGER <i>On entry:</i> whether or not the Jacobian being used has been calculated at the beginning of the current iteration. If the Jacobian has been updated then IFLAG = 1; otherwise IFLAG = 2. The Jacobian is only calculated when convergence to the current argument values has been slow.	<i>Input</i>
3:	IFAIL1 – INTEGER <i>On entry:</i> if $-6 \leq ISTATE \leq -1$, IFAIL1 specifies the IFAIL error number that would be produced were control returned to you. IFAIL1 is unspecified for values of ISTATE outside this range.	<i>Input</i>
4:	P(M) – REAL (KIND=nag_wp) array <i>On entry:</i> the current estimate of the i th argument p_i , for $i = 1, 2, \dots, m$.	<i>Input</i>
5:	M – INTEGER <i>On entry:</i> m , the number of arguments.	<i>Input</i>
6:	F(M) – REAL (KIND=nag_wp) array <i>On entry:</i> r , the residual corresponding to the current argument values, provided $1 \leq ISTATE \leq 5$ or $ISTATE = -7$. F is unspecified for other values of ISTATE.	<i>Input</i>
7:	PNORM – REAL (KIND=nag_wp) <i>On entry:</i> a quantity against which all reductions in norm are currently measured.	<i>Input</i>
8:	PNORM1 – REAL (KIND=nag_wp) <i>On entry:</i> p , the norm of the current arguments. It is set for $1 \leq ISTATE \leq 5$ and is undefined for other values of ISTATE.	<i>Input</i>
9:	EPS – REAL (KIND=nag_wp) <i>On entry:</i> gives some indication of the convergence rate. It is the current singular value modification factor (see Gay (1976)). It is zero initially and whenever convergence is	<i>Input</i>

proceeding steadily. EPS is $\epsilon^{3/8}$ or greater (where ϵ may in most cases be considered **machine precision**) when the singular values of J are approximately zero or when convergence is not being achieved. The larger the value of EPS the worse the convergence rate. When EPS becomes too large the Newton iteration is terminated.

10: D(M) – REAL (KIND=nag_wp) array *Input*

On entry: J , the singular values of the current modified Jacobian matrix. If $D(m)$ is small relative to $D(1)$ for a number of Jacobians corresponding to different argument values then the computed results should be viewed with suspicion. It could be that the matching equations do not depend significantly on some argument (which could be due to a programming error in FCN, BC, RANGE or EQN). Alternatively, the system of differential equations may be very ill-conditioned when viewed as an initial value problem, in which case D02SAF is unsuitable. This may also be indicated by some singular values being very large. These values of $D(i)$, for $i = 1, 2, \dots, m$, should not be changed.

MONIT must either be a module subprogram USED by, or declared as EXTERNAL in, the (sub) program from which D02SAF is called. Arguments denoted as *Input* must **not** be changed by this procedure.

20: PRSOL – SUBROUTINE, supplied by the NAG Library or the user. *External Procedure*

PRSOL can be used to obtain values of the solution y at a selected point z by integration across the final range $[X(1), X(NPOINT)]$. If no output is required D02HBW can be used as the actual argument.

The specification of PRSOL is:

SUBROUTINE PRSOL (Z, Y, N)

INTEGER N
REAL (KIND=nag_wp) Z, Y(N)

1: Z – REAL (KIND=nag_wp) *Input/Output*

On entry: contains x_1 on the first call. On subsequent calls Z contains its previous output value.

On exit: the next point at which output is required. The new point must be nearer $X(NPOINT)$ than the old.

If Z is set to a point outside $[X(1), X(NPOINT)]$ the process stops and control returns from D02SAF to the (sub)program from which D02SAF is called. Otherwise the next call to PRSOL is made by D02SAF at the point Z, with solution values y_1, y_2, \dots, y_n at Z contained in Y. If Z is set to $X(NPOINT)$ exactly, the final call to PRSOL is made with y_1, y_2, \dots, y_n as values of the solution at $X(NPOINT)$ produced by the integration. In general the solution values obtained at $X(NPOINT)$ from PRSOL will differ from the values obtained at this point by a call to BC. The difference between the two solutions is the residual r . You are reminded that the points $X(1), X(2), \dots, X(NPOINT)$ are available in the locations $SWP(1, 4), SWP(2, 4), \dots, SWP(NPOINT, 4)$ at all times.

2: Y(N) – REAL (KIND=nag_wp) array *Input*

On entry: the solution value y_i , for $i = 1, 2, \dots, n$, at z .

3: N – INTEGER *Input*

On entry: n , the total number of differential equations.

PRSOL must either be a module subprogram USED by, or declared as EXTERNAL in, the (sub) program from which D02SAF is called. Arguments denoted as *Input* must **not** be changed by this procedure.

- 21: W(LDW, SDW) – REAL (KIND=nag_wp) array *Output*
On exit: in the case of an error exit of the type where the point of failure is returned in SWP(1,5), the solution at this point of failure is returned in $W(i, 1)$, for $i = 1, 2, \dots, n$.
 Otherwise W is used for workspace.
- 22: LDW – INTEGER *Input*
On entry: the first dimension of the array W as declared in the (sub)program from which D02SAF is called.
Constraint: $LDW \geq \max(N, M)$.
- 23: SDW – INTEGER *Input*
On entry: the second dimension of the array W as declared in the (sub)program from which D02SAF is called.
Constraint: $SDW \geq 3 \times M + 12 + \max(11, M)$.
- 24: IFAIL – INTEGER *Input/Output*
On entry: IFAIL must be set to 0, -1 or 1. If you are unfamiliar with this argument you should refer to Section 3.4 in How to Use the NAG Library and its Documentation for details.
 For environments where it might be inappropriate to halt program execution when an error is detected, the value -1 or 1 is recommended. If the output of error messages is undesirable, then the value 1 is recommended. Otherwise, if you are not familiar with this argument, the recommended value is 0. **When the value -1 or 1 is used it is essential to test the value of IFAIL on exit.**
On exit: IFAIL = 0 unless the routine detects an error or a warning has been flagged (see Section 6).

6 Error Indicators and Warnings

If on entry IFAIL = 0 or -1, explanatory error messages are output on the current error message unit (as defined by X04AAF).

Errors or warnings detected by the routine:

IFAIL = 1

One or more of the arguments N, N1, M, LDSWP, NPOINT, ICOUNT, LDW, SDW, E, PE or YMAX is incorrectly set.

IFAIL = 2

The constraints have been violated by the initial arguments.

IFAIL = 3

The condition $X(1) < X(2) < \dots < X(NPOINT)$ (or $X(1) > X(2) > \dots > X(NPOINT)$) has been violated on a call to RANGE with the initial arguments.

IFAIL = 4

In the integration from $X(1)$ to $X(NPOINT)$ with the initial or the final arguments, the step size was reduced too far for the integration to proceed. Consider reversing the order of the points $X(1), X(2), \dots, X(NPOINT)$. If this error exit still results, it is likely that D02SAF is not a suitable method for solving the problem, or the initial choice of arguments is very poor, or the accuracy requirement specified by $E(i)$, for $i = 1, 2, \dots, n$, is too stringent.

IFAIL = 5

In the integration from $X(1)$ to $X(NPOINT)$ with the initial or final arguments, an initial step could not be found to start the integration on one of the intervals $X(i)$ to $X(i+1)$. Consider reversing the order of the points. If this error exit still results it is likely that D02SAF is not a suitable routine for solving the problem, or the initial choice of arguments is very poor, or the accuracy requirement specified by $E(i)$, for $i = 1, 2, \dots, n$, is much too stringent.

IFAIL = 6

In the integration from $X(1)$ to $X(NPOINT)$ with the initial or final arguments, the solution exceeded YMAX in magnitude (when $YMAX > 0.0$). It is likely that the initial choice of arguments was very poor or YMAX was incorrectly set.

Note: on an error with IFAIL = 4, 5 or 6 with the initial arguments, the interval in which failure occurs is contained in SWP(NPOINT, 1). If a MONIT similar to the one in Section 10 is being used then it is a simple matter to distinguish between errors using the initial and final arguments. None of the error exits IFAIL = 4, 5 or 6 should occur on the **final** integration (when computing the solution) as this integration has already been performed previously with exactly the same arguments p_i , for $i = 1, 2, \dots, m$. Seek expert help if this error occurs.

IFAIL = 7

On calculating the initial approximation to the Jacobian, the constraints were violated.

IFAIL = 8

On perturbing the arguments when calculating the initial approximation to the Jacobian, the condition $X(1) < X(2) < \dots < X(NPOINT)$ (or $X(1) > X(2) > \dots > X(NPOINT)$) is violated.

IFAIL = 9

On calculating the initial approximation to the Jacobian, the integration step size was reduced too far to make further progress (see IFAIL = 4).

IFAIL = 10

On calculating the initial approximation to the Jacobian, the initial integration step size on some interval was too small (see IFAIL = 5).

IFAIL = 11

On calculating the initial approximation to the Jacobian, the solution of the system of differential equations exceeded YMAX in magnitude (when $YMAX > 0.0$).

Note: all the error exits IFAIL = 7, 8, 9, 10 and 11 can be treated by reducing the size of some or all the elements of DP.

IFAIL = 12

On calculating the initial approximation to the Jacobian, a column of the Jacobian is found to be insignificant. This could be due to an element $DP(i)$ being too small (but nonzero) or the solution having no dependence on one of the arguments (a programming error).

Note: on an error exit with IFAIL = 7, 8, 9, 10, 11 or 12, if a perturbation of the argument p_i is the cause of the error then SWP(NPOINT, 1) will contain the floating-point value of i .

IFAIL = 13

After calculating the initial approximation to the Jacobian, the calculation of its singular value decomposition failed. It is likely that the error will never occur as it is usually associated with the Jacobian having multiple singular values. To remedy the error it should only be necessary to change the initial arguments. If the error persists it is likely that the problem has not been correctly formulated.

IFAIL = 14

The Newton iteration has failed to converge after exercising all its options. You are strongly recommended to monitor the progress of the iteration via MONIT. There are many possible reasons for the iteration not converging. Amongst the most likely are:

- (a) there is no solution;
- (b) the initial arguments are too far away from the correct arguments;
- (c) the problem is too ill-conditioned as an initial value problem for Newton's method to choose suitable corrections;
- (d) the accuracy requirements for convergence are too restrictive, that is some of the components of PE (and maybe PF) are too small – in this case the final value of this norm output via MONIT will usually be very small; or
- (e) the initial arguments are so close to the solution arguments p that the Newton iteration cannot find improved arguments. The norm output by MONIT should be very small.

IFAIL = 15

The number of iterations permitted by ICOUNT has been exceeded (in the case when ICOUNT > 0 on entry).

IFAIL = 16

IFAIL = 17

IFAIL = 18

IFAIL = 19

These indicate that there has been a serious error in an internal call. Check all subroutine calls and array dimensions. Seek expert help.

IFAIL = -99

An unexpected error has been triggered by this routine. Please contact NAG.

See Section 3.9 in How to Use the NAG Library and its Documentation for further information.

IFAIL = -399

Your licence key may have expired or may not have been installed correctly.

See Section 3.8 in How to Use the NAG Library and its Documentation for further information.

IFAIL = -999

Dynamic memory allocation failed.

See Section 3.7 in How to Use the NAG Library and its Documentation for further information.

7 Accuracy

If the iteration converges, the accuracy to which the unknown arguments are determined is usually close to that specified by you. The accuracy of the solution (output via PRSOL) depends on the error tolerances $E(i)$, for $i = 1, 2, \dots, n$. You are strongly recommended to vary all tolerances to check the accuracy of the arguments p and the solution y .

8 Parallelism and Performance

D02SAF is not thread safe and should not be called from a multithreaded user program. Please see Section 3.12.1 in How to Use the NAG Library and its Documentation for more information on thread safety.

D02SAF makes calls to BLAS and/or LAPACK routines, which may be threaded within the vendor library used by this implementation. Consult the documentation for the vendor library for further information.

Please consult the X06 Chapter Introduction for information on how to control and interrogate the OpenMP environment used within this routine. Please also consult the Users' Note for your implementation for any additional implementation-specific information.

9 Further Comments

The time taken by D02SAF depends on the complexity of the system of differential equations and on the number of iterations required. In practice, the integration of the differential system (1) is usually by far the most costly process involved. The computing time for integrating the differential equations can sometimes depend critically on the quality of the initial estimates for the arguments p . If it seems that too much computing time is required and, in particular, if the values of the residuals (output in MONIT) are much larger than expected given your knowledge of the expected solution, then the coding of FCN, EQN, RANGE and BC should be checked for errors. If no errors can be found then an independent attempt should be made to improve the initial estimates p .

In the case of an error exit in the integration of the differential system indicated by IFAIL = 4, 5, 9 or 10 you are strongly recommended to perform trial integrations with D02PFF to determine the effects of changes of the local error tolerances and of changes to the initial choice of the arguments p_i , for $i = 1, 2, \dots, m$, (that is the initial choice of p).

It is possible that by following the advice given in Section 6 an error exit with IFAIL = 7, 8, 9, 10 or 11 might be followed by one with IFAIL = 12 (or vice-versa) where the advice given is the opposite. If you are unable to refine the choice of $DP(i)$, for $i = 1, 2, \dots, n$, such that both these types of exits are avoided then the problem should be rescaled if possible or the method must be abandoned.

The choice of the 'floor' values $PF(i)$, for $i = 1, 2, \dots, m$, may be critical in the convergence of the Newton iteration. For each value i , the initial choice of p_i and the choice of $PF(i)$ should not both be very small unless it is expected that the final argument p_i will be very small and that it should be determined accurately in a **relative** sense.

For many problems it is critical that a good initial estimate be found for the arguments p or the iteration will not converge or may even break down with an error exit. There are many mathematical techniques which obtain good initial estimates for p in simple cases but which may fail to produce useful estimates in harder cases. If no such technique is available it is recommended that you try a continuation (homotopy) technique preferably based on a physical argument (e.g., the Reynolds or Prandtl number is often a suitable continuation argument). In a continuation method a sequence of problems is solved, one for each choice of the continuation argument, starting with the problem of interest. At each stage the arguments p calculated at earlier stages are used to compute a good initial estimate for the arguments at the current stage (see Hall and Watt (1976) for more details).

10 Example

This example intends to illustrate the use of the break-point and equation solving facilities of D02SAF. Most of the facilities which are common to D02SAF and D02HBF are illustrated in the example in the specification of D02HBF (which should also be consulted).

The program solves a projectile problem in two media determining the position of change of media, p_3 , and the gravity and viscosity in the second medium (p_2 represents gravity and p_4 represents viscosity).

10.1 Program Text

```
!   D02SAF Example Program Text
!   Mark 26 Release. NAG Copyright 2016.

Module d02safe_mod

!       D02SAF Example Program Module:
!       Parameters and User-defined Routines
```

```

!      .. Use Statements ..
      Use nag_library, Only: nag_wp
!      .. Implicit None Statement ..
      Implicit None
!      .. Accessibility Statements ..
      Private
      Public                                :: bc, constr, eqn, fcn, prsol, range
!      .. Parameters ..
      Real (Kind=nag_wp), Parameter        :: alpha = 0.032_nag_wp
      Real (Kind=nag_wp), Parameter        :: beta = 0.02_nag_wp
      Real (Kind=nag_wp), Parameter        :: xend = 5.0_nag_wp
      Integer, Parameter, Public           :: iset = 1, m = 4, n = 3, nin = 5,      &
                                          nout = 6
Contains
      Subroutine eqn(e,q,p,m)

!      .. Scalar Arguments ..
      Integer, Intent (In)                 :: m, q
!      .. Array Arguments ..
      Real (Kind=nag_wp), Intent (Out) :: e(q)
      Real (Kind=nag_wp), Intent (In)  :: p(m)
!      .. Executable Statements ..
      e(1) = 0.02_nag_wp - p(4) - 1.0E-5_nag_wp*p(3)
      Return
End Subroutine eqn
      Subroutine fcn(x,y,f,n,p,m,i)

!      .. Scalar Arguments ..
      Real (Kind=nag_wp), Intent (In) :: x
      Integer, Intent (In)           :: i, m, n
!      .. Array Arguments ..
      Real (Kind=nag_wp), Intent (Out) :: f(n)
      Real (Kind=nag_wp), Intent (In)  :: p(m), y(n)
!      .. Intrinsic Procedures ..
      Intrinsic                        :: cos, tan
!      .. Executable Statements ..
      f(1) = tan(y(3))
      If (i==1) Then
         f(2) = -alpha*tan(y(3))/y(2) - beta*y(2)/cos(y(3))
         f(3) = -alpha/y(2)**2
      Else
         f(2) = -p(2)*tan(y(3))/y(2) - p(4)*y(2)/cos(y(3))
         f(3) = -p(2)/y(2)**2
      End If
      Return
End Subroutine fcn
      Subroutine bc(g1,g2,p,m,n)

!      .. Scalar Arguments ..
      Integer, Intent (In)           :: m, n
!      .. Array Arguments ..
      Real (Kind=nag_wp), Intent (Out) :: g1(n), g2(n)
      Real (Kind=nag_wp), Intent (In)  :: p(m)
!      .. Executable Statements ..
      g1(1) = 0.0_nag_wp
      g1(2) = 0.5_nag_wp
      g1(3) = p(1)
      g2(1) = 0.0_nag_wp
      g2(2) = 0.45_nag_wp
      g2(3) = -1.2_nag_wp
      Return
End Subroutine bc
      Subroutine range(x,npoint,p,m)

!      .. Scalar Arguments ..
      Integer, Intent (In)           :: m, npoint
!      .. Array Arguments ..
      Real (Kind=nag_wp), Intent (In) :: p(m)
      Real (Kind=nag_wp), Intent (Out) :: x(npoint)
!      .. Executable Statements ..
      x(1) = 0.0_nag_wp

```

```

        x(2) = p(3)
        x(3) = xend
        Return
End Subroutine range
Subroutine prsol(z,y,n)

!      .. Scalar Arguments ..
      Real (Kind=nag_wp), Intent (Inout) :: z
      Integer, Intent (In) :: n
!      .. Array Arguments ..
      Real (Kind=nag_wp), Intent (In) :: y(n)
!      .. Local Scalars ..
      Integer :: i
!      .. Intrinsic Procedures ..
      Intrinsic :: abs
!      .. Executable Statements ..
      If (z==0.0E0_nag_wp) Then
        Write (nout,*)
        Write (nout,*) '      Z      Y(1)      Y(2)      Y(3)'
      End If
      Write (nout,99999) z, (y(i),i=1,n)
      z = z + 0.5_nag_wp
      If (abs(z-xend)<0.25_nag_wp) Then
        z = xend
      End If
      Return

99999  Format (1X,F9.3,3F10.4)
End Subroutine prsol
Function constr(p,m)

!      .. Function Return Value ..
      Logical :: constr
!      .. Scalar Arguments ..
      Integer, Intent (In) :: m
!      .. Array Arguments ..
      Real (Kind=nag_wp), Intent (In) :: p(m)
!      .. Intrinsic Procedures ..
      Intrinsic :: any
!      .. Executable Statements ..
      If (any(p(1:m)<0.0_nag_wp) .Or. p(3)>5.0_nag_wp) Then
        constr = .False.
      Else
        constr = .True.
      End If
      Return
End Function constr
End Module d02safe_mod
Program d02safe

!      D02SAF Example Main Program

!      .. Use Statements ..
      Use nag_library, Only: d02saf, d02sas, nag_wp, x04abf
      Use d02safe_mod, Only: bc, constr, eqn, fcn, iset, m, n, nin, nout, &
        prsol, range
!      .. Implicit None Statement ..
      Implicit None
!      .. Local Scalars ..
      Real (Kind=nag_wp) :: ymax
      Integer :: i, icount, ifail, ldswp, ldw, nl, &
        npoint, outchn, sdw
!      .. Local Arrays ..
      Real (Kind=nag_wp), Allocatable :: dp(:), e(:), p(:), pe(:), pf(:), &
        swp(:,:), w(:,:)
!      .. Intrinsic Procedures ..
      Intrinsic :: max
!      .. Executable Statements ..
      Write (nout,*) 'D02SAF Example Program Results'
!      Skip heading in data file
      Read (nin,*)

```



```

Read (nin,*) npoint
n1 = n
sdw = 3*m + 23
ldswp = npoint
ldw = max(m,n)
Allocate (dp(m),e(n),p(m),pe(m),pf(m),swp(ldswp,6),w(ldw,sdw))

outchn = nout
Read (nin,*) icount
Read (nin,*) ymax
Read (nin,*) pe(1:m)
Read (nin,*) pf(1:m)
Read (nin,*) dp(1:m)
Read (nin,*) e(1:n)
Call x04abf(iset,outchn)
swp(1:npoint-1,1:3) = 0.0_nag_wp
Read (nin,*) p(1:m)

!      * To obtain monitoring information, replace the name d02sas by d02hbx
!      in the next statement and USE nag_library : d02hbx

!      ifail: behaviour on error exit
!      =0 for hard exit, =1 for quiet-soft, =-1 for noisy-soft
ifail = 1
Call d02saf(p,m,n,n1,pe,pf,e,dp,npoint,swp,ldswp,icount,range,bc,fcn,      &
    eqn,constr,ymax,d02sas,prsol,w,ldw,sdw,ifail)

If (ifail/=0) Then
    Write (nout,99999) ifail
End If
If (ifail>=4 .And. ifail<=12) Then
    Write (nout,99998) 'SWP(NPOINT,1) = ', swp(npoint,1)
    If (ifail<=6) Then
        Write (nout,99998) 'SWP(1,5) = ', swp(1,5)
        Write (nout,*) '      i      W(i,1) '
        Write (nout,99997)(i,w(i,1),i=1,n)
    End If
End If

99999 Format (1X,/,1X,' ** D02SAF returned with IFAIL = ',I5)
99998 Format (1X,A,F10.4)
99997 Format (1X,I4,1X,E10.3)
End Program d02safe

```

10.2 Program Data

D02SAF Example Program Data

3	:	npoint
0	:	icount
0.0	:	ymax
1.0E-3 1.0E-3 1.0E-3 1.0E-3	:	pe(1:m)
1.0E-6 1.0E-6 1.0E-6 1.0E-6	:	pf(1:m)
0.0 0.0 0.0 0.0	:	dp(1:m)
1.0E-5 1.0E-5 1.0E-5	:	e (1:n)
1.2 0.032 2.5 0.02	:	p (1:m)

10.3 Program Results

D02SAF Example Program Results

Z	Y(1)	Y(2)	Y(3)
0.000	0.0000	0.5000	1.1753
0.500	1.0881	0.4127	1.0977
1.000	1.9501	0.3310	0.9802
1.500	2.5768	0.2582	0.7918
2.000	2.9606	0.2019	0.4796
2.500	3.0958	0.1773	0.0245

3.000	2.9861	0.1935	-0.4353
3.500	2.6289	0.2409	-0.7679
4.000	2.0181	0.3047	-0.9767
4.500	1.1454	0.3759	-1.1099
5.000	0.0000	0.4500	-1.2000

