

## NAG Toolbox

### nag\_fit\_1dspline\_auto (e02be)

#### 1 Purpose

nag\_fit\_1dspline\_auto (e02be) computes a cubic spline approximation to an arbitrary set of data points. The knots of the spline are located automatically, but a single argument must be specified to control the trade-off between closeness of fit and smoothness of fit.

#### 2 Syntax

```
[n, lamda, c, fp, wrk, iwrk, ifail] = nag_fit_1dspline_auto(start, x, y, w, s, n,
lamda, wrk, iwrk, 'm', m, 'nest', nest)

[n, lamda, c, fp, wrk, iwrk, ifail] = e02be(start, x, y, w, s, n, lamda, wrk,
iwrk, 'm', m, 'nest', nest)
```

**Note:** the interface to this routine has changed since earlier releases of the toolbox:

At Mark 22: *lwrk* was removed from the interface.

#### 3 Description

nag\_fit\_1dspline\_auto (e02be) determines a smooth cubic spline approximation  $s(x)$  to the set of data points  $(x_r, y_r)$ , with weights  $w_r$ , for  $r = 1, 2, \dots, m$ .

The spline is given in the B-spline representation

$$s(x) = \sum_{i=1}^{n-4} c_i N_i(x), \quad (1)$$

where  $N_i(x)$  denotes the normalized cubic B-spline defined upon the knots  $\lambda_i, \lambda_{i+1}, \dots, \lambda_{i+4}$ .

The total number  $n$  of these knots and their values  $\lambda_1, \dots, \lambda_n$  are chosen automatically by the function. The knots  $\lambda_5, \dots, \lambda_{n-4}$  are the interior knots; they divide the approximation interval  $[x_1, x_m]$  into  $n - 7$  sub-intervals. The coefficients  $c_1, c_2, \dots, c_{n-4}$  are then determined as the solution of the following constrained minimization problem:

minimize

$$\eta = \sum_{i=5}^{n-4} \delta_i^2 \quad (2)$$

subject to the constraint

$$\theta = \sum_{r=1}^m \epsilon_r^2 \leq S, \quad (3)$$

where  $\delta_i$  stands for the discontinuity jump in the third order derivative of  $s(x)$  at the interior knot  $\lambda_i$ ,  
 $\epsilon_r$  denotes the weighted residual  $w_r(y_r - s(x_r))$ ,  
 and  $S$  is a non-negative number to be specified by you.

The quantity  $\eta$  can be seen as a measure of the (lack of) smoothness of  $s(x)$ , while closeness of fit is measured through  $\theta$ . By means of the argument  $S$ , ‘the smoothing factor’, you can then control the balance between these two (usually conflicting) properties. If  $S$  is too large, the spline will be too smooth and signal will be lost (underfit); if  $S$  is too small, the spline will pick up too much noise (overfit). In the extreme cases the function will return an interpolating spline ( $\theta = 0$ ) if  $S$  is set to zero,

and the weighted least squares cubic polynomial ( $\eta = 0$ ) if  $S$  is set very large. Experimenting with  $S$  values between these two extremes should result in a good compromise. (See Section 9.2 for advice on choice of  $S$ .)

The method employed is outlined in Section 9.3 and fully described in Dierckx (1975), Dierckx (1981) and Dierckx (1982). It involves an adaptive strategy for locating the knots of the cubic spline (depending on the function underlying the data and on the value of  $S$ ), and an iterative method for solving the constrained minimization problem once the knots have been determined.

Values of the computed spline, or of its derivatives or definite integral, can subsequently be computed by calling `nag_fit_1dspline_eval` (e02bb), `nag_fit_1dspline_deriv` (e02bc) or `nag_fit_1dspline_integ` (e02bd), as described in Section 9.4.

## 4 References

Dierckx P (1975) An algorithm for smoothing, differentiating and integration of experimental data using spline functions *J. Comput. Appl. Math.* **1** 165–184

Dierckx P (1981) An improved algorithm for curve fitting with spline functions *Report TW54* Department of Computer Science, Katholieke Universiteit Leuven

Dierckx P (1982) A fast algorithm for smoothing data on a rectangular grid while using spline functions *SIAM J. Numer. Anal.* **19** 1286–1304

Reinsch C H (1967) Smoothing by spline functions *Numer. Math.* **10** 177–183

## 5 Parameters

### 5.1 Compulsory Input Parameters

1: **start** – CHARACTER(1)

Must be set to 'C' or 'W'.

**start** = 'C'

The function will build up the knot set starting with no interior knots. No values need be assigned to the arguments **n**, **lamda**, **wrk** or **iwrk**.

**start** = 'W'

The function will restart the knot-placing strategy using the knots found in a previous call of the function. In this case, the arguments **n**, **lamda**, **wrk**, and **iwrk** must be unchanged from that previous call. This warm start can save much time in searching for a satisfactory value of **s**.

*Constraint:* **start** = 'C' or 'W'.

2: **x(m)** – REAL (KIND=nag\_wp) array

The values  $x_r$  of the independent variable (abscissa)  $x$ , for  $r = 1, 2, \dots, m$ .

*Constraint:*  $x_1 < x_2 < \dots < x_m$ .

3: **y(m)** – REAL (KIND=nag\_wp) array

The values  $y_r$  of the dependent variable (ordinate)  $y$ , for  $r = 1, 2, \dots, m$ .

4: **w(m)** – REAL (KIND=nag\_wp) array

The values  $w_r$  of the weights, for  $r = 1, 2, \dots, m$ . For advice on the choice of weights, see Section 2.1.2 in the E02 Chapter Introduction.

*Constraint:* **w**( $r$ ) > 0.0, for  $r = 1, 2, \dots, m$ .

- 5: **s** – REAL (KIND=nag\_wp)  
 The smoothing factor,  $S$ .  
 If  $S = 0.0$ , the function returns an interpolating spline.  
 If  $S$  is smaller than *machine precision*, it is assumed equal to zero.  
 For advice on the choice of  $S$ , see Section 3 and Section 9.2.  
*Constraint:*  $s \geq 0.0$ .
- 6: **n** – INTEGER  
 If the warm start option is used, the value of **n** must be left unchanged from the previous call.
- 7: **lamda(nest)** – REAL (KIND=nag\_wp) array  
 If the warm start option is used, the values **lamda**(1), **lamda**(2), ..., **lamda**(**n**) must be left unchanged from the previous call.
- 8: **wrk(lwrk)** – REAL (KIND=nag\_wp) array  
*lwrk*, the dimension of the array, must satisfy the constraint  $lwrk \geq 4 \times \mathbf{m} + 16 \times \mathbf{nest} + 41$ .  
 If the warm start option is used on entry, the values **wrk**(1), ..., **wrk**(**n**) must be left unchanged from the previous call.
- 9: **iwrk(nest)** – INTEGER array  
 If the warm start option is used, on entry, the values **iwrk**(1), ..., **iwrk**(**n**) must be left unchanged from the previous call.  
 This array is used as workspace.

## 5.2 Optional Input Parameters

- 1: **m** – INTEGER  
*Default:* the dimension of the arrays **x**, **y**, **w**. (An error is raised if these dimensions are not equal.)  
*m*, the number of data points.  
*Constraint:*  $\mathbf{m} \geq 4$ .
- 2: **nest** – INTEGER  
*Default:* the dimension of the arrays **lamda**, **iwrk**. (An error is raised if these dimensions are not equal.)  
 An overestimate for the number,  $n$ , of knots required.  
*Constraint:*  $\mathbf{nest} \geq 8$ . In most practical situations,  $\mathbf{nest} = \mathbf{m}/2$  is sufficient. **nest** never needs to be larger than  $\mathbf{m} + 4$ , the number of knots needed for interpolation ( $s = 0.0$ ).

## 5.3 Output Parameters

- 1: **n** – INTEGER  
 The total number,  $n$ , of knots of the computed spline.
- 2: **lamda(nest)** – REAL (KIND=nag\_wp) array  
 The knots of the spline, i.e., the positions of the interior knots **lamda**(5), **lamda**(6), ..., **lamda**(**n** - 4) as well as the positions of the additional knots  

$$\mathbf{lamda}(1) = \mathbf{lamda}(2) = \mathbf{lamda}(3) = \mathbf{lamda}(4) = x_1$$

and

$$\mathbf{lamda}(n-3) = \mathbf{lamda}(n-2) = \mathbf{lamda}(n-1) = \mathbf{lamda}(n) = x_m$$

needed for the B-spline representation.

3: **c**(**nest**) – REAL (KIND=nag\_wp) array

The coefficient  $c_i$  of the B-spline  $N_i(x)$  in the spline approximation  $s(x)$ , for  $i = 1, 2, \dots, n-4$ .

4: **fp** – REAL (KIND=nag\_wp)

The sum of the squared weighted residuals,  $\theta$ , of the computed spline approximation. If **fp** = 0.0, this is an interpolating spline. **fp** should equal **s** within a relative tolerance of 0.001 unless  $n = 8$  when the spline has no interior knots and so is simply a cubic polynomial. For knots to be inserted, **s** must be set to a value below the value of **fp** produced in this case.

5: **wrk**(*lwrk*) – REAL (KIND=nag\_wp) array

6: **iwrk**(**nest**) – INTEGER array

This array is used as workspace.

7: **ifail** – INTEGER

**ifail** = 0 unless the function detects an error (see Section 5).

## 6 Error Indicators and Warnings

Errors or warnings detected by the function:

**ifail** = 1

On entry, **start**  $\neq$  'C' or 'W',  
 or **m** < 4,  
 or **s** < 0.0,  
 or **s** = 0.0 and **nest** < **m** + 4,  
 or **nest** < 8,  
 or *lwrk* <  $4 \times \mathbf{m} + 16 \times \mathbf{nest} + 41$ .

**ifail** = 2

The weights are not all strictly positive.

**ifail** = 3

The values of  $\mathbf{x}(r)$ , for  $r = 1, 2, \dots, \mathbf{m}$ , are not in strictly increasing order.

**ifail** = 4

The number of knots required is greater than **nest**. Try increasing **nest** and, if necessary, supplying larger arrays for the arguments **lamda**, **c**, **wrk** and **iwrk**. However, if **nest** is already large, say **nest** > **m**/2, then this error exit may indicate that **s** is too small.

**ifail** = 5

The iterative process used to compute the coefficients of the approximating spline has failed to converge. This error exit may occur if **s** has been set very small. If the error persists with increased **s**, contact NAG.

**ifail** = -99

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

**ifail** = –399

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

**ifail** = –999

Dynamic memory allocation failed.

If **ifail** = 4 or 5, a spline approximation is returned, but it fails to satisfy the fitting criterion (see (2) and (3)) – perhaps by only a small amount, however.

## 7 Accuracy

On successful exit, the approximation returned is such that its weighted sum of squared residuals  $\theta$  (as in (3)) is equal to the smoothing factor  $S$ , up to a specified relative tolerance of 0.001 – except that if  $n = 8$ ,  $\theta$  may be significantly less than  $S$ : in this case the computed spline is simply a weighted least squares polynomial approximation of degree 3, i.e., a spline with no interior knots.

## 8 Further Comments

### 8.1 Timing

The time taken for a call of `nag_fit_1dspline_auto` (e02be) depends on the complexity of the shape of the data, the value of the smoothing factor  $S$ , and the number of data points. If `nag_fit_1dspline_auto` (e02be) is to be called for different values of  $S$ , much time can be saved by setting **start** = 'W' after the first call.

### 8.2 Choice of $S$

If the weights have been correctly chosen (see Section 2.1.2 in the E02 Chapter Introduction), the standard deviation of  $w_r y_r$  would be the same for all  $r$ , equal to  $\sigma$ , say. In this case, choosing the smoothing factor  $S$  in the range  $\sigma^2(m \pm \sqrt{2m})$ , as suggested by Reinsch (1967), is likely to give a good start in the search for a satisfactory value. Otherwise, experimenting with different values of  $S$  will be required from the start, taking account of the remarks in Section 3.

In that case, in view of computation time and memory requirements, it is recommended to start with a very large value for  $S$  and so determine the least squares cubic polynomial; the value returned in **fp**, call it  $\theta_0$ , gives an upper bound for  $S$ . Then progressively decrease the value of  $S$  to obtain closer fits – say by a factor of 10 in the beginning, i.e.,  $S = \theta_0/10$ ,  $S = \theta_0/100$ , and so on, and more carefully as the approximation shows more details.

The number of knots of the spline returned, and their location, generally depend on the value of  $S$  and on the behaviour of the function underlying the data. However, if `nag_fit_1dspline_auto` (e02be) is called with **start** = 'W', the knots returned may also depend on the smoothing factors of the previous calls. Therefore if, after a number of trials with different values of  $S$  and **start** = 'W', a fit can finally be accepted as satisfactory, it may be worthwhile to call `nag_fit_1dspline_auto` (e02be) once more with the selected value for  $S$  but now using **start** = 'C'. Often, `nag_fit_1dspline_auto` (e02be) then returns an approximation with the same quality of fit but with fewer knots, which is therefore better if data reduction is also important.

### 8.3 Outline of Method Used

If  $S = 0$ , the requisite number of knots is known in advance, i.e.,  $n = m + 4$ ; the interior knots are located immediately as  $\lambda_i = x_{i-2}$ , for  $i = 5, 6, \dots, n - 4$ . The corresponding least squares spline (see `nag_fit_1dspline_knots` (e02ba)) is then an interpolating spline and therefore a solution of the problem.

If  $S > 0$ , a suitable knot set is built up in stages (starting with no interior knots in the case of a cold start but with the knot set found in a previous call if a warm start is chosen). At each stage, a spline is fitted to the data by least squares (see `nag_fit_1dspline_knots` (e02ba)) and  $\theta$ , the weighted sum of squares of residuals, is computed. If  $\theta > S$ , new knots are added to the knot set to reduce  $\theta$  at the next stage. The new knots are located in intervals where the fit is particularly poor, their number depending on the value of  $S$  and on the progress made so far in reducing  $\theta$ . Sooner or later, we find that  $\theta \leq S$  and

at that point the knot set is accepted. The function then goes on to compute the (unique) spline which has this knot set and which satisfies the full fitting criterion specified by (2) and (3). The theoretical solution has  $\theta = S$ . The function computes the spline by an iterative scheme which is ended when  $\theta = S$  within a relative tolerance of 0.001. The main part of each iteration consists of a linear least squares computation of special form, done in a similarly stable and efficient manner as in `nag_fit_1dspline_knots` (e02ba).

An exception occurs when the function finds at the start that, even with no interior knots ( $n = 8$ ), the least squares spline already has its weighted sum of squares of residuals  $\leq S$ . In this case, since this spline (which is simply a cubic polynomial) also has an optimal value for the smoothness measure  $\eta$ , namely zero, it is returned at once as the (trivial) solution. It will usually mean that  $S$  has been chosen too large.

For further details of the algorithm and its use, see Dierckx (1981).

## 8.4 Evaluation of Computed Spline

The value of the computed spline at a given value  $\mathbf{x}$  may be obtained in the double variable  $\mathbf{s}$  by the call:

```
[s, ifail] = e02bb(lamda, c, x);
```

where  $\mathbf{n}$ ,  $\mathbf{lamda}$  and  $\mathbf{c}$  are the output arguments of `nag_fit_1dspline_auto` (e02be).

The values of the spline and its first three derivatives at a given value  $\mathbf{x}$  may be obtained in the double array  $\mathbf{s}$  of dimension at least 4 by the call:

```
[s, ifail] = e02bc(lamda, c, x, left);
```

where if  $\mathbf{left} = 1$ , left-hand derivatives are computed and if  $\mathbf{left} \neq 1$ , right-hand derivatives are calculated. The value of  $\mathbf{left}$  is only relevant if  $\mathbf{x}$  is an interior knot (see `nag_fit_1dspline_deriv` (e02bc)).

The value of the definite integral of the spline over the interval  $\mathbf{x}(1)$  to  $\mathbf{x}(\mathbf{m})$  can be obtained in the double variable  $\mathbf{dint}$  by the call:

```
[dint, ifail] = e02bd(lamda, c);
```

(see `nag_fit_1dspline_integ` (e02bd)).

## 9 Example

This example reads in a set of data values, followed by a set of values of  $\mathbf{s}$ . For each value of  $\mathbf{s}$  it calls `nag_fit_1dspline_auto` (e02be) to compute a spline approximation, and prints the values of the knots and the B-spline coefficients  $c_i$ .

The program includes code to evaluate the computed splines, by calls to `nag_fit_1dspline_eval` (e02bb), at the points  $x_r$  and at points mid-way between them. These values are not printed out, however; instead the results are illustrated by plots of the computed splines, together with the data points (indicated by  $\times$ ) and the positions of the knots (indicated by vertical lines): the effect of decreasing  $\mathbf{s}$  can be clearly seen.

### 9.1 Program Text

```
function e02be_example

fprintf('e02be example results\n\n');

start = 'C';
data = [0.0  -1.100  1.0;
        0.5  -0.372  2.0;
        1.0   0.431  1.5;
        1.5   1.690  1.0;
        2.0   2.110  3.0;
        2.5   3.100  1.0;
        3.0   4.230  0.5;
        4.0   4.350  1.0;
```

```

        4.5   4.810  2.0;
        5.0   4.610  2.5;
        5.5   4.790  1.0;
        6.0   5.230  3.0;
        7.0   6.350  1.0;
        7.5   7.190  2.0;
        8.0   7.970  1.0];
x = data(:,1);
y = data(:,2);
w = data(:,3);
m = size(data,1);

s = [1 0.5 0.1];

nest = m + 4;
n     = nag_int(nest);
lamda = zeros(nest,1);
wrk   = zeros(4*m+16*nest+41, 1);
iwrk  = zeros(nest, 1, nag_int_name);
start = 'Cold';

xs(1:2:2*m-1) = x;
xs(2:2:2*m-2) = (x(1:m-1)+x(2:m))/2;

for is = 1:3
    % Get spline
    [n, lamda, c, fp, wrk, iwrk, ifail] = ...
    e02be( ...
        start, x, y, w, s(is), n, lamda, wrk, iwrk);

    % Print details of spline
    fprintf('\nCalling with smoothing factor S = %5.2f\n\n', s(is));
    fprintf('      Knots      Coeffs\n');
    fprintf('%4d%20.4f\n', 1, c(1));
    for j = 2:n-5
        fprintf('%4d%10.4f%10.4f\n', j, lamda(j+2), c(j));
    end
    fprintf('%4d%20.4f\n\n', n-4, c(n-4));
    fprintf('Weighted sum of squared residuals = %7.4f\n', fp);
    if fp==0
        fprintf('(The spline is an interpolating spline)\n');
    elseif n==8
        fprintf('(The spline is the weighted least-squares cubic polynomial)\n');
    end
    fprintf('\n');

    % Evaluate at x and mid-points (xs)
    for i = 1:2*m-1
        [fit(i,is), ifail] = e02bb( ...
            lamda, c, xs(i), 'ncap7', n);
    end

    start = 'Warm';
end

fig1 = figure;
hold on
plot(x,y,'*','Color','Green')
plot(xs,fit(:,1));
xlabel('x');
title('Evaluation of B-spline representation, S = 1.0');
legend('Evaluation points','B-spline','Location','NorthWest');
hold off;

fig2 = figure;
hold on
plot(x,y,'*','Color','Green')
plot(xs,fit(:,2));
xlabel('x');
title('Evaluation of B-spline representation, S = 0.5');
legend('Evaluation points','B-spline','Location','NorthWest');

```

```

hold off;

fig3 = figure;
hold on
plot(x,y,'*','Color','Green')
plot(xs,fit(:,3));
xlabel('x');
title('Evaluation of B-spline representation, S = 0.1');
legend('Evaluation points','B-spline','Location','NorthWest');
hold off;

```

## 9.2 Program Results

e02be example results

Calling with smoothing factor S = 1.00

	Knots	Coeffs
1		-1.3201
2	0.0000	1.3542
3	4.0000	5.5510
4	8.0000	4.7031
5		8.2277

Weighted sum of squared residuals = 1.0003

Calling with smoothing factor S = 0.50

	Knots	Coeffs
1		-1.1072
2	0.0000	-0.6571
3	1.0000	0.4350
4	2.0000	2.8061
5	4.0000	4.6824
6	5.0000	4.6416
7	6.0000	5.1976
8	8.0000	6.9008
9		7.9979

Weighted sum of squared residuals = 0.5001

Calling with smoothing factor S = 0.10

	Knots	Coeffs
1		-1.0901
2	0.0000	-0.6401
3	1.0000	0.0334
4	1.5000	1.6390
5	2.0000	2.1243
6	3.0000	4.5591
7	4.0000	4.2174
8	4.5000	4.9105
9	5.0000	4.5475
10	5.5000	4.6960
11	6.0000	5.7370
12	8.0000	6.8179
13		7.9953

Weighted sum of squared residuals = 0.1000







