

NAG Fortran Library Routine Document

E02BEF

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

E02BEF computes a cubic spline approximation to an arbitrary set of data points. The knots of the spline are located automatically, but a single parameter must be specified to control the trade-off between closeness of fit and smoothness of fit.

2 Specification

```

SUBROUTINE E02BEF (START, M, X, Y, W, S, NEST, N, LAMDA, C, FP, WRK,
1 LWRK, IWRK, IFAIL)
INTEGER M, NEST, N, LWRK, IWRK(NEST), IFAIL
double precision X(M), Y(M), W(M), S, LAMDA(NEST), C(NEST), FP,
1 WRK(LWRK)
CHARACTER*1 START

```

3 Description

E02BEF 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 routine. 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 δ_i stands for the discontinuity jump in the third order derivative of $s(x)$ at the interior knot λ_i ,
 ϵ_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 η can be seen as a measure of the (lack of) smoothness of $s(x)$, while closeness of fit is measured through θ . By means of the parameter S , 'the smoothing factor', you will 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 routine 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 8.2 for advice on choice of S .)

The method employed is outlined in Section 8.3 and fully described in Dierckx (1975), Dierckx (1981a) 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 E02BBF, E02BCF or E02BDF, as described in Section 8.4.

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 (1981a) 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

- 1: START – CHARACTER*1 *Input*
On entry: must be set to 'C' or 'W'.
 START = 'C' (Cold start)
 The routine will build up the knot set starting with no interior knots. No values need be assigned to the parameters N, LAMDA, WRK or IWRK.
 START = 'W' (Warm start)
 The routine will restart the knot-placing strategy using the knots found in a previous call of the routine. In this case, the parameters 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: M – INTEGER *Input*
On entry: m , the number of data points.
Constraint: $M \geq 4$.
- 3: X(M) – *double precision* array *Input*
On entry: the values x_r of the independent variable (abscissa) x , for $r = 1, 2, \dots, m$.
Constraint: $x_1 < x_2 < \dots < x_m$.
- 4: Y(M) – *double precision* array *Input*
On entry: the values y_r of the dependent variable (ordinate) y , for $r = 1, 2, \dots, m$.
- 5: W(M) – *double precision* array *Input*
On entry: 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: $W(r) > 0$, for $r = 1, 2, \dots, m$.
- 6: S – *double precision* *Input*
On entry: the smoothing factor, S .

If $S = 0.0$, the routine 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 Sections 3 and 8.2.

Constraint: $S \geq 0.0$.

- 7: NEST – INTEGER *Input*
On entry: an over-estimate for the number, n , of knots required.
Constraint: $NEST \geq 8$. In most practical situations, $NEST = M/2$ is sufficient. NEST never needs to be larger than $M + 4$, the number of knots needed for interpolation ($S = 0.0$).
- 8: N – INTEGER *Input/Output*
On entry: if the warm start option is used, the value of N must be left unchanged from the previous call.
On exit: the total number, n , of knots of the computed spline.
- 9: LAMDA(NEST) – *double precision* array *Input/Output*
On entry: if the warm start option is used, the values $LAMDA(1), LAMDA(2), \dots, LAMDA(N)$ must be left unchanged from the previous call.
On exit: the knots of the spline, i.e., the positions of the interior knots $LAMDA(5), LAMDA(6), \dots, LAMDA(N - 4)$ as well as the positions of the additional knots
$$LAMDA(1) = LAMDA(2) = LAMDA(3) = LAMDA(4) = x_1$$
and
$$LAMDA(N - 3) = LAMDA(N - 2) = LAMDA(N - 1) = LAMDA(N) = x_m$$
needed for the B-spline representation.
- 10: C(NEST) – *double precision* array *Output*
On exit: the coefficient c_i of the B-spline $N_i(x)$ in the spline approximation $s(x)$, for $i = 1, 2, \dots, n - 4$.
- 11: FP – *double precision* *Output*
On exit: the sum of the squared weighted residuals, θ , 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.
- 12: WRK(LWRK) – *double precision* array *Communication Array*
If the warm start option is used on entry, the values $WRK(1), \dots, WRK(n)$ must be left unchanged from the previous call.
- 13: LWRK – INTEGER *Input*
On entry: the dimension of the array WRK as declared in the (sub)program from which E02BEF is called.
Constraint: $LWRK \geq 4 \times M + 16 \times NEST + 41$.
- 14: IWRK(NEST) – INTEGER array *Communication Array*
If the warm start option is used, on entry, the values $IWRK(1), \dots, IWRK(n)$ must be left unchanged from the previous call.
This array is used as workspace.

15: IFAIL – INTEGER

Input/Output

On entry: IFAIL must be set to 0, -1 or 1. If you are unfamiliar with this parameter you should refer to Chapter P01 for details.

On exit: IFAIL = 0 unless the routine detects an error (see Section 6).

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 parameter the recommended value is 0. **When the value -1 or 1 is used it is essential to test the value of IFAIL on exit.**

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

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 $LWRK < 4 \times M + 16 \times NEST + 41$.

IFAIL = 2

The weights are not all strictly positive.

IFAIL = 3

The values of $X(r)$, for $r = 1, 2, \dots, M$, 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 parameters 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, consult NAG.

If IFAIL = 4 or 5, a spline approximation is returned, but it fails to satisfy the fitting criterion (see (2) and (3) in Section 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 FP is equal to the smoothing factor S, up to a specified relative tolerance of 0.001 – except that if $n = 8$, FP 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 E02BEF depends on the complexity of the shape of the data, the value of the smoothing factor S , and the number of data points. If E02BEF 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 σ , 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 for FP , call it FP_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., $S = FP_0/10$, $S = FP_0/100$, 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 E02BEF 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 E02BEF once more with the selected value for S but now using $START = 'C'$. Often, E02BEF 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 E02BAF) 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 E02BAF) and θ , the weighted sum of squares of residuals, is computed. If $\theta > S$, new knots are added to the knot set to reduce θ 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 θ . Sooner or later, we find that $\theta \leq S$ and at that point the knot set is accepted. The routine 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 routine 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 E02BAF.

An exception occurs when the routine 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 η , 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 (1981a).

8.4 Evaluation of Computed Spline

The value of the computed spline at a given value X may be obtained in the *double precision* variable S by the call:

```
CALL E02BBF(N,LAMDA,C,X,S,IFAIL)
```

where N , $LAMDA$ and C are the output parameters of E02BEF.

The values of the spline and its first three derivatives at a given value X may be obtained in the **double precision** array S of dimension at least 4 by the call:

```
CALL E02BCF(N,LAMDA,C,X,LEFT,S,IFAIL)
```

where if $LEFT = 1$, left-hand derivatives are computed and if $LEFT \neq 1$, right-hand derivatives are calculated. The value of $LEFT$ is only relevant if X is an interior knot (see E02BCF).

The value of the definite integral of the spline over the interval $X(1)$ to $X(M)$ can be obtained in the **double precision** variable $DINT$ by the call:

```
CALL E02BDF(N,LAMDA,C,DINT,IFAIL)
```

(see E02BDF).

9 Example

This example reads in a set of data values, followed by a set of values of S . For each value of S it calls E02BEF 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 E02BBF, 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 S can be clearly seen. (The plots were obtained by calling NAG Graphical Supplement routine J06FAF.)

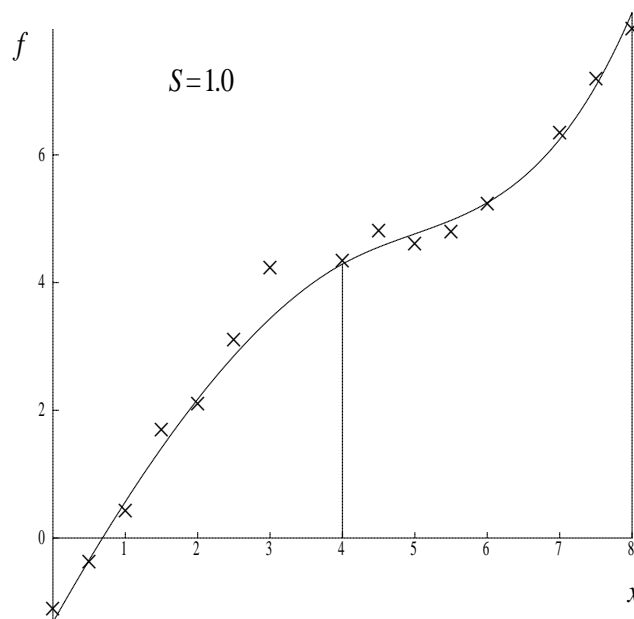


Figure 1

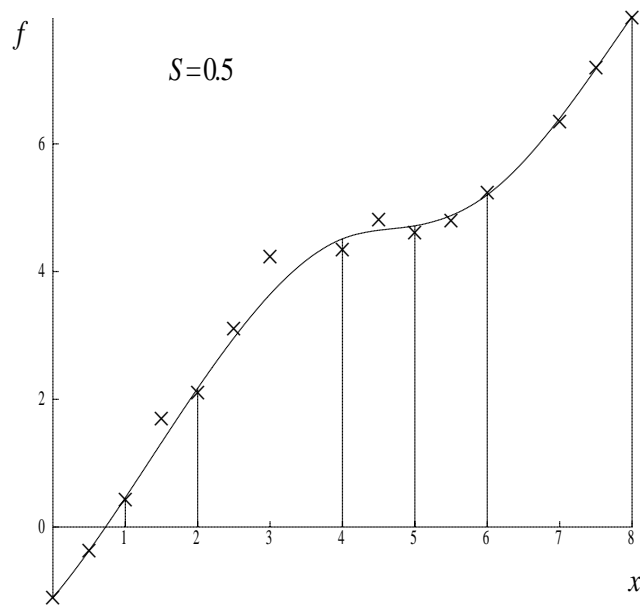


Figure 2

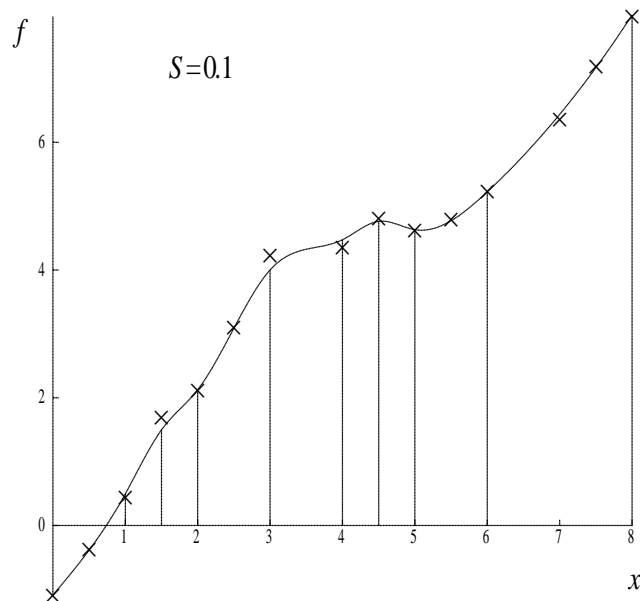


Figure 3

9.1 Program Text

```

*      E02BEF Example Program Text
*      Mark 20 Revised. NAG Copyright 2001.
*      .. Parameters ..
      INTEGER          NIN, NOUT
      PARAMETER       (NIN=5,NOUT=6)
      INTEGER          MMAX, NEST, LWRK
      PARAMETER       (MMAX=50,NEST=MMAX+4,LWRK=4*MMAX+16*NEST+41)
*      .. Local Scalars ..
      DOUBLE PRECISION FP, S, TXR
      INTEGER          IFAIL, J, M, N, R
      CHARACTER        START
*      .. Local Arrays ..
      DOUBLE PRECISION C(NEST), K(NEST), SP(2*MMAX-1), W(MMAX),

```

```

+          WRK(LWRK), X(MMAX), Y(MMAX)
INTEGER      IWRK(NEST)
*
* .. External Subroutines ..
EXTERNAL     E02BBF, E02BEF
*
* .. Executable Statements ..
WRITE (NOUT,*) 'E02BEF Example Program Results'
*
* Skip heading in data file
READ (NIN,*)
*
* Input the number of data points, followed by the data points (X),
* the function values (Y) and the weights (W).
READ (NIN,*) M
IF (M.GT.0 .AND. M.LE.MMAX) THEN
  DO 20 R = 1, M
    READ (NIN,*) X(R), Y(R), W(R)
20  CONTINUE
  START = 'C'
*
* Read in successive values of S until end of data file.
40  READ (NIN,*,END=120) S
*
* Determine the spline approximation.
IFAIL = 0
*
*
* CALL E02BEF(START,M,X,Y,W,S,NEST,N,K,C,FP,WRK,LWRK,IWRK,IFAIL)
*
* Evaluate the spline at each X point and midway between
* X points, saving the results in SP.
DO 60 R = 1, M
  IFAIL = 0
*
* CALL E02BBF(N,K,C,X(R),SP((R-1)*2+1),IFAIL)
*
60  CONTINUE
DO 80 R = 1, M - 1
  IFAIL = 0
  TXR = (X(R)+X(R+1))/2
*
* CALL E02BBF(N,K,C,TXR,SP(R*2),IFAIL)
*
80  CONTINUE
*
* Output the results.
WRITE (NOUT,*)
WRITE (NOUT,99999) 'Calling with smoothing factor S =', S
WRITE (NOUT,*)
WRITE (NOUT,*)
+          '          B-Spline'
+          '          J          Knot K(J+2)   Coefficient C(J)'
WRITE (NOUT,99998) 1, C(1)
DO 100 J = 2, N - 5
  WRITE (NOUT,99997) J, K(J+2), C(J)
100 CONTINUE
WRITE (NOUT,99998) N - 4, C(N-4)
WRITE (NOUT,*)
WRITE (NOUT,99999) 'Weighted sum of squared residuals FP =', FP
IF (FP.EQ.0.0D0) THEN
  WRITE (NOUT,*) '(The spline is an interpolating spline)'
ELSE IF (N.EQ.8) THEN
  WRITE (NOUT,*)
+          '(The spline is the weighted least-squares cubic polynomial)'
END IF
WRITE (NOUT,*)
START = 'W'
GO TO 40
END IF
120 STOP
*
99999 FORMAT (1X,A,1P,E12.3)
99998 FORMAT (11X,I4,16X,F16.4)
99997 FORMAT (11X,I4,2F16.4)
END

```


9.2 Program Data

E02BEF Example Program Data

```

15
0.0000E+00 -1.1000E+00 1.00 M, the number of data points
5.0000E-01 -3.7200E-01 2.00 X, Y, W, abscissa, ordinate and weight
1.0000E+00 4.3100E-01 1.50
1.5000E+00 1.6900E+00 1.00
2.0000E+00 2.1100E+00 3.00
2.5000E+00 3.1000E+00 1.00
3.0000E+00 4.2300E+00 0.50
4.0000E+00 4.3500E+00 1.00
4.5000E+00 4.8100E+00 2.00
5.0000E+00 4.6100E+00 2.50
5.5000E+00 4.7900E+00 1.00
6.0000E+00 5.2300E+00 3.00
7.0000E+00 6.3500E+00 1.00
7.5000E+00 7.1900E+00 2.00
8.0000E+00 7.9700E+00 1.00 End of data points
1.0 S, smoothing factor
0.5 S, smoothing factor
0.1 S, smoothing factor

```

9.3 Program Results

E02BEF Example Program Results

Calling with smoothing factor S = 1.000E+00

J	Knot K(J+2)	B-Spline Coefficient C(J)
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 FP = 1.000E+00

Calling with smoothing factor S = 5.000E-01

J	Knot K(J+2)	B-Spline Coefficient C(J)
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 FP = 5.001E-01

Calling with smoothing factor S = 1.000E-01

J	Knot K(J+2)	B-Spline Coefficient C(J)
1		-1.0900
2	0.0000	-0.6422
3	1.0000	0.0369
4	1.5000	1.6353
5	2.0000	2.1274
6	3.0000	4.5526
7	4.0000	4.2225
8	4.5000	4.9108
9	5.0000	4.4159

10	6.0000	5.4794
11	8.0000	6.8308
12		7.9935

Weighted sum of squared residuals FP = 1.000E-01
