

NAG Toolbox for MATLAB

e02dc

1 Purpose

e02dc computes a bicubic spline approximation to a set of data values, given on a rectangular grid in the x - y plane. 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 Syntax

```
[nx, lamda, ny, mu, c, fp, wrk, iwrk, ifail] = e02dc(start, x, y, f, s,
nx, lamda, ny, mu, wrk, iwrk, 'mx', mx, 'my', my, 'nxest', nxest,
'nyest', nyest, 'lwrk', lwrk, 'liwrk', liwrk)
```

3 Description

e02dc determines a smooth bicubic spline approximation $s(x, y)$ to the set of data points $(x_q, y_r, f_{q,r})$, for $q = 1, 2, \dots, m_x$ and $r = 1, 2, \dots, m_y$.

The spline is given in the B-spline representation

$$s(x, y) = \sum_{i=1}^{n_x-4} \sum_{j=1}^{n_y-4} c_{ij} M_i(x) N_j(y), \quad (1)$$

where $M_i(x)$ and $N_j(y)$ denote normalized cubic B-splines, the former defined on the knots λ_i to λ_{i+4} and the latter on the knots μ_j to μ_{j+4} . For further details, see Hayes and Halliday 1974 for bicubic splines and de Boor 1972 for normalized B-splines.

The total numbers n_x and n_y of these knots and their values $\lambda_1, \dots, \lambda_{n_x}$ and μ_1, \dots, μ_{n_y} are chosen automatically by the function. The knots $\lambda_5, \dots, \lambda_{n_x-4}$ and $\mu_5, \dots, \mu_{n_y-4}$ are the interior knots; they divide the approximation domain $[x_1, x_{m_x}] \times [y_1, y_{m_y}]$ into $(n_x - 7) \times (n_y - 7)$ subpanels $[\lambda_i, \lambda_{i+1}] \times [\mu_j, \mu_{j+1}]$, for $i = 4, 5, \dots, n_x - 4$, $j = 4, 5, \dots, n_y - 4$. Then, much as in the curve case (see e02be), the coefficients c_{ij} are determined as the solution of the following constrained minimization problem:

minimize

$$\eta, \quad (2)$$

subject to the constraint

$$\theta = \sum_{q=1}^{m_x} \sum_{r=1}^{m_y} \epsilon_{q,r}^2 \leq S, \quad (3)$$

where η is a measure of the (lack of) smoothness of $s(x, y)$. Its value depends on the discontinuity jumps in $s(x, y)$ across the boundaries of the subpanels. It is zero only when there are no discontinuities and is positive otherwise, increasing with the size of the jumps (see Dierckx 1982 for details).

$\epsilon_{q,r}$ denotes the residual $f_{q,r} - s(x_q, y_r)$,
and S is a nonnegative number to be specified by you.

By means of the parameter **s**, ‘the smoothing factor’, you will then control the balance between smoothness and closeness of fit, as measured by the sum of squares of residuals in (3). 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 least-squares bicubic 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 8.3 for advice on choice of s .)

The method employed is outlined in Section 8.5 and fully described in Dierckx 1982 and Dierckx 1981a. It involves an adaptive strategy for locating the knots of the bicubic 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 can subsequently be computed by calling e02de or e02df as described in Section 8.6.

4 References

de Boor C 1972 On calculating with B-splines *J. Approx. Theory* **6** 50–62

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

Hayes J G and Halliday J 1974 The least-squares fitting of cubic spline surfaces to general data sets *J. Inst. Math. Appl.* **14** 89–103

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

5 Parameters

5.1 Compulsory Input Parameters

1: **start** – string

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

start = 'C' (Cold start)

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

start = 'W' (Warm start)

The function will restart the knot-placing strategy using the knots found in a previous call of the function. In this case, the parameters **nx**, **ny**, **lamda**, **mu**, **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(mx)** – double array

$x(q)$ must be set to x_q , the x co-ordinate of the q th grid point along the x axis, for $q = 1, 2, \dots, m_x$.

Constraint: $x_1 < x_2 < \dots < x_{m_x}$.

3: **y(my)** – double array

$y(r)$ must be set to y_r , the y co-ordinate of the r th grid point along the y axis, for $r = 1, 2, \dots, m_y$.

Constraint: $y_1 < y_2 < \dots < y_{m_y}$.

4: **f(mx × my)** – double array

$f(m_y \times (q - 1) + r)$ must contain the data value $f_{q,r}$, for $q = 1, 2, \dots, m_x$ and $r = 1, 2, \dots, m_y$.

5: **s – double scalar**

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 Sections 3 and 8.3.

Constraint: $s \geq 0.0$.

6: **nx – int32 scalar**

If the warm start option is used, the value of **nx** must be left unchanged from the previous call.

7: **lamda(nxest) – double array**

If the warm start option is used, the values **lamda**(1), **lamda**(2), ..., **lamda**(**nx**) must be left unchanged from the previous call.

8: **ny – int32 scalar**

If the warm start option is used, the value of **ny** must be left unchanged from the previous call.

9: **mu(nyest) – double array**

If the warm start option is used, the values **mu**(1), **mu**(2), ..., **mu**(**ny**) must be left unchanged from the previous call.

10: **wrk(lwrk) – double array**

If the warm start option is used, on entry, the values **wrk**(1), ..., **wrk**(4) must be left unchanged from the previous call.

This array is used as workspace.

11: **iwrk(liwrk) – int32 array**

If the warm start option is used, on entry, the values **iwrk**(1), ..., **iwrk**(3) must be left unchanged from the previous call.

This array is used as workspace.

5.2 Optional Input Parameters1: **mx – int32 scalar**

Default: The dimension of the array **x**.

m_x , the number of grid points along the x axis.

Constraint: $mx \geq 4$.

2: **my – int32 scalar**

Default: The dimension of the array **y**.

m_y , the number of grid points along the y axis.

Constraint: $my \geq 4$.

3: **nxest – int32 scalar**4: **nyest – int32 scalar**

Default: For **nxest**, the dimension of the array **lamda**. For **nyest**, the dimension of the array **mu**.
an upper bound for the number of knots n_x and n_y required in the x - and y -directions respectively.

In most practical situations, $\mathbf{nxest} = m_x/2$ and $\mathbf{nyest} = m_y/2$ is sufficient. \mathbf{nxest} and \mathbf{nyest} never need to be larger than $m_x + 4$ and $m_y + 4$ respectively, the numbers of knots needed for interpolation ($\mathbf{s} = 0.0$). See also Section 8.4.

Constraint: $\mathbf{nxest} \geq 8$ and $\mathbf{nyest} \geq 8$.

5: **lwrk – int32 scalar**

Default: The dimension of the array **wrk**.

Constraint: $\mathbf{lwrk} \geq 4 \times (\mathbf{mx} + \mathbf{my}) + 11 \times (\mathbf{nxest} + \mathbf{nyest}) + \mathbf{nxest} \times \mathbf{my} + \max(\mathbf{my}, \mathbf{nxest}) + 54$.

6: **liwrk – int32 scalar**

Default: The dimension of the array **iwrk**.

Constraint: $\mathbf{liwrk} \geq 3 + \mathbf{mx} + \mathbf{my} + \mathbf{nxest} + \mathbf{nyest}$.

5.3 Input Parameters Omitted from the MATLAB Interface

None.

5.4 Output Parameters

1: **nx – int32 scalar**

The total number of knots, n_x , of the computed spline with respect to the x variable.

2: **lamda(nxest) – double array**

Contains the complete set of knots λ_i associated with the x variable, i.e., the interior knots **lamda(5)**, **lamda(6)**, ..., **lamda(nx – 4)** as well as the additional knots

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

and

$$\mathbf{lamda}(\mathbf{nx} - 3) = \mathbf{lamda}(\mathbf{nx} - 2) = \mathbf{lamda}(\mathbf{nx} - 1) = \mathbf{lamda}(\mathbf{nx}) = \mathbf{x}(\mathbf{mx})$$

needed for the B-spline representation.

3: **ny – int32 scalar**

The total number of knots, n_y , of the computed spline with respect to the y variable.

4: **mu(nyest) – double array**

Contains the complete set of knots μ_i associated with the y variable, i.e., the interior knots **mu(5)**, **mu(6)**, ..., **mu(ny – 4)** as well as the additional knots

$$\mathbf{mu}(1) = \mathbf{mu}(2) = \mathbf{mu}(3) = \mathbf{mu}(4) = \mathbf{y}(1)$$

and

$$\mathbf{mu}(\mathbf{ny} - 3) = \mathbf{mu}(\mathbf{ny} - 2) = \mathbf{mu}(\mathbf{ny} - 1) = \mathbf{mu}(\mathbf{ny}) = \mathbf{y}(\mathbf{my})$$

needed for the B-spline representation.

5: **c((nxest – 4) × (nyest – 4)) – double array**

The coefficients of the spline approximation. $\mathbf{c}((n_y - 4) \times (i - 1) + j)$ is the coefficient c_{ij} defined in Section 3.

6: **fp – double scalar**

The sum of squared 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 **nx** = **ny** = 8, when the spline has no interior knots and so is simply a bicubic polynomial. For knots to be inserted, **s** must be set to a value below the value of **fp** produced in this case.

7: **wrk(lwrk) – double array**

If the warm start option is used, on entry, the values **wrk**(1), ..., **wrk**(4) must be left unchanged from the previous call.

This array is used as workspace.

8: **iwrk(liwrk) – int32 array**

If the warm start option is used, on entry, the values **iwrk**(1), ..., **iwrk**(3) must be left unchanged from the previous call.

This array is used as workspace.

9: **ifail – int32 scalar**

0 unless the function detects an error (see Section 6).

6 Error Indicators and Warnings

Errors or warnings detected by the function:

ifail = 1

On entry, **start** \neq 'C' or 'W',

or **mx** < 4,

or **my** < 4,

or **s** < 0.0,

or **s** = 0.0 and **nxest** < **mx** + 4,

or **s** = 0.0 and **nyest** < **my** + 4,

or **nxest** < 8,

or **nyest** < 8,

or **lwrk** < $4 \times (\mathbf{mx} + \mathbf{my}) + 11 \times (\mathbf{nxest} + \mathbf{nyest}) + \mathbf{nxest} \times \mathbf{my} + \max(\mathbf{my}, \mathbf{nxest}) + 54$,

or **liwrk** < $3 + \mathbf{mx} + \mathbf{my} + \mathbf{nxest} + \mathbf{nyest}$.

ifail = 2

The values of **x**(*q*), for *q* = 1, 2, ..., **mx**, are not in strictly increasing order.

ifail = 3

The values of **y**(*r*), for *r* = 1, 2, ..., **my**, are not in strictly increasing order.

ifail = 4

The number of knots required is greater than allowed by **nxest** and **nyest**. Try increasing **nxest** and/or **nyest** and, if necessary, supplying larger arrays for the parameters **lamda**, **mu**, **c**, **wrk** and **iwrk**. However, if **nxest** and **nyest** are already large, say **nxest** > **mx**/2 and **nyest** > **my**/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 sum of squared residuals **fp** is equal to the smoothing factor **s**, up to a specified relative tolerance of 0.001 – except that if $n_x = 8$ and $n_y = 8$, **fp** may be significantly less than **s**: in this case the computed spline is simply the least-squares bicubic 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 e02dc depends on the complexity of the shape of the data, the value of the smoothing factor **s**, and the number of data points. If e02dc is to be called for different values of **s**, much time can be saved by setting **start** = 'W' after the first call.

8.2 Weighting of Data Points

e02dc does not allow individual weighting of the data values. If these were determined to widely differing accuracies, it may be better to use e02dd. The computation time would be very much longer, however.

8.3 Choice of **s**

If the standard deviation of $f_{q,r}$ is the same for all q and r (the case for which e02dc is designed – see Section 8.2.) and known to be equal, at least approximately, to σ , say, then following Reinsch 1967 and choosing the smoothing factor **s** in the range $\sigma^2(m \pm \sqrt{2m})$, where $m = m_x m_y$, is likely to give a good start in the search for a satisfactory value. If the standard deviations vary, the sum of their squares over all the data points could be used. 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 bicubic 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 e02dc 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 e02dc once more with the selected value for **s** but now using **start** = 'C'. Often, e02dc then returns an approximation with the same quality of fit but with fewer knots, which is therefore better if data reduction is also important.

8.4 Choice of **nxest** and **nyest**

The number of knots may also depend on the upper bounds **nxest** and **nyest**. Indeed, if at a certain stage in e02dc the number of knots in one direction (say n_x) has reached the value of its upper bound (**nxest**), then from that moment on all subsequent knots are added in the other (y) direction. Therefore you have the option of limiting the number of knots the function locates in any direction. For example, by setting **nxest** = 8 (the lowest allowable value for **nxest**), you can indicate that he wants an approximation which is a simple cubic polynomial in the variable x .

8.5 Outline of Method Used

If **s** = 0, the requisite number of knots is known in advance, i.e., $n_x = m_x + 4$ and $n_y = m_y + 4$; the interior knots are located immediately as $\lambda_i = x_{i-2}$ and $\mu_j = y_{j-2}$, for $i = 5, 6, \dots, n_x - 4$ and $j = 5, 6, \dots, n_y - 4$.

The corresponding least-squares spline is then an interpolating spline and therefore a solution of the problem.

If $s > 0$, suitable knot sets are 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 bicubic spline is fitted to the data by least-squares, and θ , the sum of squares of residuals, is computed. If $\theta > s$, new knots are added to one knot set or the other so as 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 sets are accepted. The function then goes on to compute the (unique) spline which has these knot sets 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 e02ba for least-squares curve-fitting.

An exception occurs when the function finds at the start that, even with no interior knots ($n_x = n_y = 8$), the least-squares spline already has its sum of residuals $\leq s$. In this case, since this spline (which is simply a bicubic 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 1982.

8.6 Evaluation of Computed Spline

The values of the computed spline at the points $(\mathbf{x}(r), \mathbf{y}(r))$, for $r = 1, 2, \dots, \mathbf{m}$, may be obtained in the double array **ff** (see e02de), of length at least \mathbf{m} , by the following code:

```
[ff, ifail] = e02de(x, y, lamda, mu, c);
```

where **lamda**, **mu** and **c** have the same values as output from e02dc.

To evaluate the computed spline on a \mathbf{mx} by \mathbf{my} rectangular grid of points in the x - y plane, which is defined by the x co-ordinates stored in $\mathbf{x}(q)$, for $q = 1, 2, \dots, \mathbf{mx}$, and the y co-ordinates stored in $\mathbf{y}(r)$, for $r = 1, 2, \dots, \mathbf{my}$, returning the results in the double array **ff** (see e02df) which is of length at least $\mathbf{mx} \times \mathbf{my}$, the following call may be used:

```
[ff, ifail] = e02df(x, y, lamda, mu, c);
```

where **lamda**, **mu** and **c** have the same values as output from e02dc. The result of the spline evaluated at grid point (q, r) is returned in element $(\mathbf{my} \times (q - 1) + r)$ of the array **ff**.

9 Example

```
start = 'C';
x = [0;
     0.5;
     1;
     1.5;
     2;
     2.5;
     3;
     3.5;
     4;
     4.5;
     5];
y = [0;
     0.5;
     1;
     1.5;
     2;
     2.5;
     3;
     3.5;
     4];
```

```
f = [1;  
      0.88758;  
      0.5403;  
      0.070736999999999999;  
      -0.41515;  
      -0.80114;  
      -0.97998999999999999;  
      -0.93446;  
      -0.65664;  
      1.5;  
      1.3564;  
      0.820449999999999999;  
      0.10611;  
      -0.62422;  
      -1.2317;  
      -1.485;  
      -1.3047;  
      -0.98547;  
      2.06;  
      1.7552;  
      1.0806;  
      0.15147;  
      -0.832290000000000001;  
      -1.6023;  
      -1.97;  
      -1.8729;  
      -1.4073;  
      2.57;  
      2.124;  
      1.3508;  
      0.17684;  
      -1.0404;  
      -2.0029;  
      -2.475;  
      -2.3511;  
      -1.6741;  
      3;  
      2.6427;  
      1.6309;  
      0.21221;  
      -1.2484;  
      -2.2034;  
      -2.97;  
      -2.8094;  
      -1.9809;  
      3.5;  
      3.1715;  
      1.8611;  
      0.24458;  
      -1.4565;  
      -2.864;  
      -3.265;  
      -3.2776;  
      -2.2878;  
      4.04;  
      3.5103;  
      2.0612;  
      0.28595;  
      -1.6946;  
      -3.2046;  
      -3.96;  
      -3.7958;  
      -2.6146;  
      4.5;  
      3.9391;  
      2.4314;  
      0.31632;  
      -1.8627;  
      -3.6351;  
      -4.455;
```



```

-4.2141;
-2.9314;
5.04;
4.3879;
2.7515;
0.35369;
-2.0707;
-4.0057;
-4.97;
-4.6823;
-3.2382;
5.505;
4.8367;
2.9717;
0.38505;
-2.2888;
-4.4033;
-5.445;
-5.1405;
-3.595;
6;
5.2755;
3.2418;
0.42442;
-2.4769;
-4.8169;
-5.93;
-5.6387;
-3.9319];
s = 0.1;
nx = int32(0);
lamda = zeros(15, 1);
ny = int32(-1077703852);
mu = zeros(13, 1);
wrk = zeros(592, 1);
iwrk = zeros(51, 1, 'int32');
[nxOut, lamdaOut, nyOut, muOut, c, fp, wrkOut, iwrkOut, ifail] = ...
    e02dc(start, x, y, f, s, nx, lamda, ny, mu, wrk, iwrk)

```

```

nxOut =
    10
lamdaOut =
    0
    0
    0
    0
    1.5000
    2.5000
    5.0000
    5.0000
    5.0000
    5.0000
    0
    0
    0
    0
    0
nyOut =
    13
muOut =
    0
    0
    0
    0
    1.0000
    2.0000
    2.5000
    3.0000
    3.5000
    4.0000

```

```
      4.0000
      4.0000
      4.0000
c =
    array elided
fp =
    0.1000
wrkOut =
    array elided
iwrkOut =
    array elided
ifail =
      0
```
