

NAG Toolbox for MATLAB**Chapter Introduction****C05 – Roots of One or More Transcendental Equations****Contents**

1	Scope of the Chapter	2
2	Background to the Problems	2
2.1	A Single Equation	2
2.2	Systems of Equations	2
3	Recommendations on Choice and Use of Available Functions	2
3.1	Zeros of Functions of One Variable	2
3.2	Solution of Sets of Nonlinear Equations	3
4	Decision Trees	4
5	Index	5
6	References	5

1 Scope of the Chapter

This chapter is concerned with the calculation of real zeros of continuous real functions of one or more variables. (Complex equations must be expressed in terms of the equivalent larger system of real equations.)

2 Background to the Problems

The chapter divides naturally into two parts.

2.1 A Single Equation

The first deals with the real zeros of a real function of a single variable $f(x)$.

There are three functions with simple calling sequences. The first assumes that you can determine an initial interval $[a, b]$ within which the desired zero lies, (that is, where $f(a) \times f(b) < 0$), and outside which all other zeros lie. The function then systematically subdivides the interval to produce a final interval containing the zero. This final interval has a length bounded by your specified error requirements; the end of the interval where the function has smallest magnitude is returned as the zero. This function is guaranteed to converge to a **simple** zero of the function. (Here we define a simple zero as a zero corresponding to a sign-change of the function; none of the available functions are capable of making any finer distinction.) However, as with the other functions described below, a non-simple zero might be determined and it is left to you to check for this. The algorithm used is due to Brent 1973.

The two other functions are both designed for the case where you are unable to specify an interval containing the simple zero. The first function starts from an initial point and performs a search for an interval containing a simple zero. If such an interval is computed then the method described above is used next to determine the zero accurately. The second method uses a ‘continuation’ method based on a secant iteration. A sequence of subproblems is solved; the first of these is trivial and the last is the actual problem of finding a zero of $f(x)$. The intermediate problems employ the solutions of earlier problems to provide initial guesses for the secant iterations used to calculate their solutions.

Three other functions are also supplied. They employ reverse communication and use the same core algorithms as the functions described above.

2.2 Systems of Equations

The functions in the second part of this chapter are designed to solve a set of nonlinear equations in n unknowns

$$f_i(x) = 0, \quad i = 1, 2, \dots, n, \quad x = (x_1, x_2, \dots, x_n)^T, \quad (1)$$

where T stands for transpose.

It is assumed that the functions are continuous and differentiable so that the matrix of first partial derivatives of the functions, the **Jacobian** matrix $J_{ij}(x) = \left(\frac{\partial f_i}{\partial x_j} \right)$ evaluated at the point x , exists, though it may not be possible to calculate it directly.

The functions f_i must be independent, otherwise there will be an infinity of solutions and the methods will fail. However, even when the functions are independent the solutions may not be unique. Since the methods are iterative, an initial guess at the solution has to be supplied, and the solution located will usually be the one closest to this initial guess.

3 Recommendations on Choice and Use of Available Functions

3.1 Zeros of Functions of One Variable

The functions can be divided into two classes. There are three functions (c05av, c05ax and c05az) all written in reverse communication form and three (c05ad, c05ag and c05aj) written in direct communication form. The direct communication functions are designed for inexperienced users and, in particular, for solving problems where the function $f(x)$ whose zero is to be calculated, can be coded as a . These

functions find the zero by using the same core algorithms as the reverse communication functions. Experienced users are recommended to use the reverse communication functions directly as they permit you more control of the calculation. Indeed, if the zero-finding process is embedded in a much larger program then the reverse communication functions should always be used.

The recommendation as to which function should be used depends mainly on whether you can supply an interval $[a, b]$ containing the zero; that is, where $f(a) \times f(b) < 0$. If the interval can be supplied, then c05ad (or, in reverse communication, c05az) should be used, in general. This recommendation should be qualified in the case when the only interval which can be supplied is very long relative to your error requirements **and** you can also supply a good approximation to the zero. In this case c05aj (or, in reverse communication, c05ax) **may** prove more efficient (though these latter functions will not provide the error bound available from c05az).

If an interval containing the zero cannot be supplied then you must choose between c05ag (or, in reverse communication, c05av followed by c05az) and c05aj (or, in reverse communication, c05ax). c05ag first determines an interval containing the zero, and then proceeds as in c05ad; it is particularly recommended when you do not have a good initial approximation to the zero. If a good initial approximation to the zero is available then c05aj is to be preferred. Since neither of these latter functions has guaranteed convergence to the zero, you are recommended to experiment with both in case of difficulty.

3.2 Solution of Sets of Nonlinear Equations

The solution of a set of nonlinear equations

$$f_i(x_1, x_2, \dots, x_n) = 0, \quad i = 1, 2, \dots, n \quad (2)$$

can be regarded as a special case of the problem of finding a minimum of a sum of squares

$$s(x) = \sum_{i=1}^m [f_i(x_1, x_2, \dots, x_n)]^2, \quad (m \geq n). \quad (3)$$

So the functions in Chapter E04 are relevant as well as the special nonlinear equations functions.

The functions for solving a set of nonlinear equations can also be divided into classes. There are four functions (c05nb, c05nc, c05pb and c05pc) all written in direct communication form and two (c05nd and c05pd) written in reverse communication form. The direct communication functions are designed for inexperienced users and, in particular, these functions require the f_i (and possibly their derivatives) to be calculated in user-supplied (sub)programs. These should be set up carefully so the Library functions can work as efficiently as possible. Experienced users are recommended to use the reverse communication functions as they permit you more control of the calculation. Indeed, if the zero-finding process is embedded in a much larger program then the reverse communication functions should always be used.

The main decision you have to make is whether to supply the derivatives $\frac{\partial f_i}{\partial x_j}$. It is advisable to do so if possible, since the results obtained by algorithms which use derivatives are generally more reliable than those obtained by algorithms which do not use derivatives.

c05pb and c05pc (or, in reverse communication, c05pd) require you to provide the derivatives, whilst c05nb and c05nc (or, in reverse communication, c05nd) do not. c05nb and c05pb are easy-to-use functions; greater flexibility may be obtained using c05nc and c05pc (or, in reverse communication, c05nd and c05pd), but these have longer parameter lists. c05za is provided for use in conjunction with c05pb and c05pc to check the user-supplied derivatives for consistency with the functions themselves. You are strongly advised to make use of this function whenever c05pb or c05pc is used.

Firstly, the calculation of the functions and their derivatives should be ordered so that **cancellation errors** are avoided. This is particularly important in a function that uses these quantities to build up estimates of higher derivatives.

Secondly, **scaling** of the variables has a considerable effect on the efficiency of a function. The problem should be designed so that the elements of x are of similar magnitude. The same comment applies to the functions, i.e., all the f_i should be of comparable size.

The accuracy is usually determined by the accuracy parameters of the functions, but the following points may be useful

- (i) Greater accuracy in the solution may be requested by choosing smaller input values for the accuracy parameters. However, if unreasonable accuracy is demanded, rounding errors may become important and cause a failure.
- (ii) Some idea of the accuracies of the x_i may be obtained by monitoring the progress of the function to see how many figures remain unchanged during the last few iterations.
- (iii) An approximation to the error in the solution x is given by e where e is the solution to the set of linear equations

$$J(x)e = -f(x)$$

where $f(x) = (f_1(x), f_2(x), \dots, f_n(x))^T$ (see Chapter F04).

Note that the QR decomposition of J is available from c05nc and c05pc (or, in reverse communication, c05nd and c05pd) so that

$$Re = -Q^T f$$

and $Q^T f$ is also provided by these functions.

- (iv) If the functions $f_i(x)$ are changed by small amounts ϵ_i , for $i = 1, 2, \dots, n$, then the corresponding change in the solution x is given approximately by σ , where σ is the solution of the set of linear equations

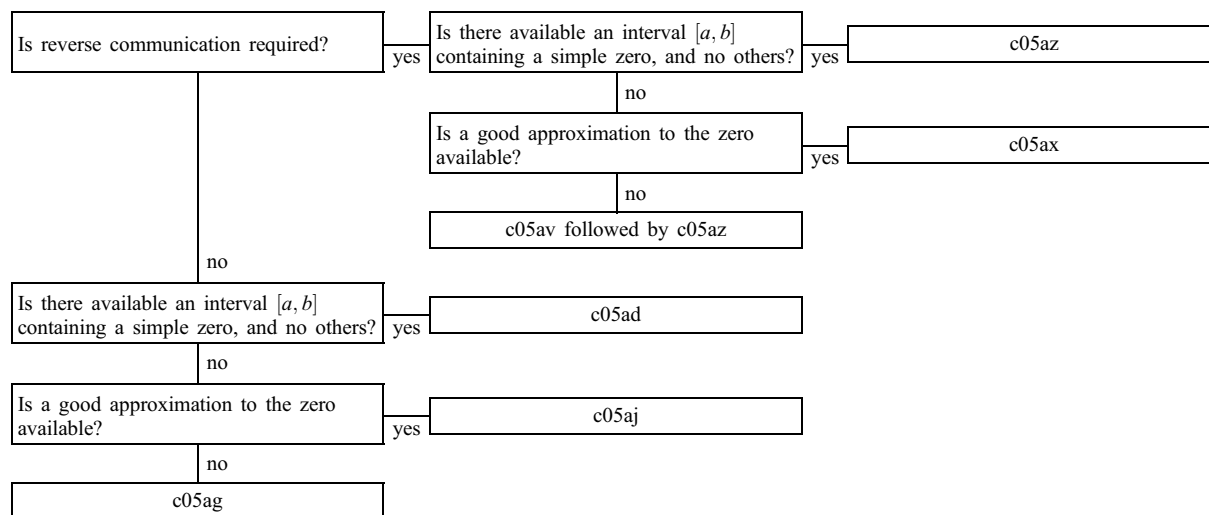
$$J(x)\sigma = -\epsilon$$

(see Chapter F04).

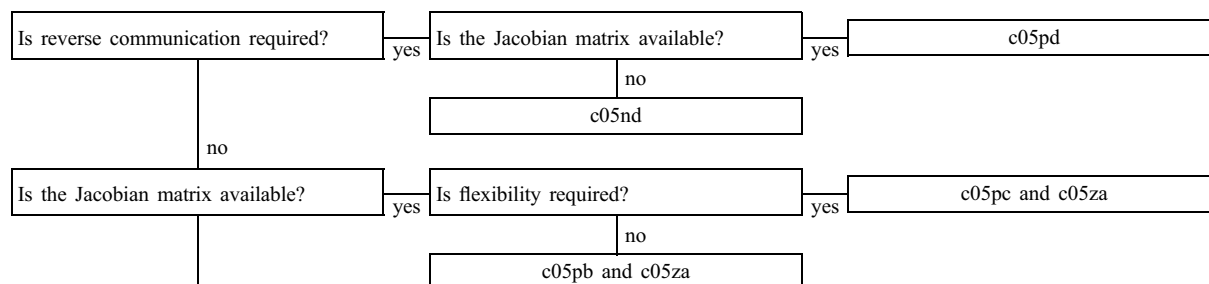
Thus one can estimate the sensitivity of x to any uncertainties in the specification of $f_i(x)$, for $i = 1, 2, \dots, n$. As noted above, the sophisticated functions c05nc and c05pc (or, in reverse communication, c05nd and c05pd) provide the QR decomposition of J .

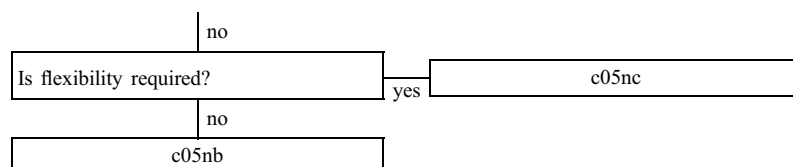
4 Decision Trees

Tree 1: Functions of One Variable



Tree 2: Functions of several variables





5 Index

Zeros of functions of one variable:

Direct communication:

binary search followed by Bus and Dekker algorithm c05ag
 Bus and Dekker algorithm c05ad
 continuation method c05aj

Reverse communication:

binary search c05av
 Bus and Dekker algorithm c05az
 continuation method c05ax

Zeros of functions of several variables:

Checking Routine:

Checks user-supplied Jacobian c05za

Direct communication:

easy-to-use,

derivatives required c05pb
 no derivatives required c05nb

sophisticated c05nc

sophisticated, derivatives required c05pc

Reverse Communication:

sophisticated c05nd

sophisticated,
 derivatives required c05pd

6 References

Brent R P 1973 *Algorithms for Minimization Without Derivatives* Prentice–Hall

Gill P E and Murray W 1976 Algorithms for the solution of the nonlinear least-squares problem *Report NAC 71* National Physical Laboratory

Moré J J, Garbow B S and Hillstom K E 1980 User guide for MINPACK-1 *Technical Report ANL-80-74* Argonne National Laboratory

Ortega J M and Rheinboldt W C 1970 *Iterative Solution of Nonlinear Equations in Several Variables* Academic Press

Rabinowitz P 1970 *Numerical Methods for Nonlinear Algebraic Equations* Gordon and Breach