

Chapter 8

Numerical Solution of Nonlinear Equations and Systems

Consider a scalar nonlinear equation:

$$F(x) = 0,$$

where $F(x)$ is a given function. Very often, equations of this type cannot be solved analytically. For example, no analytic solution is available when $F(x)$ is a high degree algebraic polynomial or when $F(x)$ is a transcendental function, such as $F(x) = \sin x - \frac{1}{2}x$ or $F(x) = e^{-x} + \cos x$. Then, a numerical method must be employed for computing the solution x approximately. Such methods are often referred to as the methods of rootfinding, because the number x that solves the equation $F(x) = 0$ is called its root.

Along with the scalar nonlinear equations, we will also consider systems of such equations:

$$\mathbf{F}(\mathbf{x}) = \mathbf{0},$$

where $\mathbf{F}(\mathbf{x})$ is a given vector-function of the vector argument $\mathbf{x} = (x_1, x_2, \dots, x_n)$. For example, if

$$\mathbf{F}(\mathbf{x}) = \begin{bmatrix} F_1(x_1, x_2) \\ F_2(x_1, x_2) \end{bmatrix} = \begin{bmatrix} x_1^2 + x_2^2 - 25 \\ x_2 - x_1^2 \end{bmatrix},$$

then the nonlinear system $\mathbf{F}(\mathbf{x}) = \mathbf{0}$ can be written in components as follows:

$$\begin{aligned} x_1^2 + x_2^2 - 25 &= 0, \\ x_2 - x_1^2 &= 0. \end{aligned}$$

When solving numerically the scalar nonlinear equations $F(x) = 0$ or systems $\mathbf{F}(\mathbf{x}) = \mathbf{0}$, one typically needs to address two issues. First, the solutions (roots) need to be isolated, i.e., the appropriate domains of the independent variable(s) need to be identified that will only contain one solution each. Then, the solutions (roots) need to be “refined,” i.e., actually computed with a prescribed accuracy.

There are no universal approaches to solving the problem of the isolation of roots. One can use graphs, look into the intervals of monotonicity of the function $F(x)$ on which it changes sign, and employ other special “tricks.” There is only one important class of functions, namely, algebraic polynomials with real coefficients, for which this problem has been solved completely in the most general form. The solution is given by the Sturm theorem, see, e.g., [Dör82, Section 24].