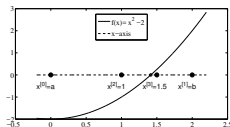


§1.1: The bisection method

Solving nonlinear equations

MA385/530 – Numerical Analysis

September 2017



Linear equations are of the form:

$$\textit{find } x \textit{ such that } ax + b = 0$$

and are easy to solve. Some nonlinear problems are also easy to solve, e.g.,

$$\textit{find } x \textit{ such that } ax^2 + bx + c = 0.$$

Cubic and quartic equations also have solutions for which we can obtain a formula. But most equations do not have simple formulae for their solutions, so numerical methods are needed.

References

- Chap. 1 of Süli and Mayers (Introduction to Numerical Analysis). We'll follow this pretty closely in lectures.
- Stewart (*Afternotes ...*), Lectures 1–5. A well-presented introduction, with lots of diagrams to give an intuitive introduction.
- Chapter 4 of Moler's "Numerical Computing with MATLAB". Gives a brief introduction to the methods we study, and description of MATLAB functions for solving these problems.
- The proof of the convergence of Newton's Method is based on the presentation in Thm 3.2 of Epperson.

Our generic problem is:

*Let f be a continuous function on the interval $[a, b]$.
Find $\tau \in [a, b]$ such that $f(\tau) = 0$.*

Here f is some specified function, and τ is the **solution** to $f(x) = 0$.

This leads to two natural questions:

- (1) How do we know there is a solution?
- (2) How do we find it?

The following gives *sufficient* conditions for the existence of a solution:

Proposition

Let f be a real-valued function that is defined and continuous on a bounded closed interval $[a, b] \subset \mathbb{R}$. Suppose that $f(a)f(b) \leq 0$. Then there exists $\tau \in [a, b]$ such that $f(\tau) = 0$.

So now we know there is a solution τ to $f(x) = 0$, but how to we actually solve it? **Usually we don't!** Instead we construct a sequence of estimates $\{x_0, x_1, x_2, x_3, \dots\}$ that **converge** to the true solution. So now we have to answer these questions:

- (1) How can we construct the sequence x_0, x_1, \dots ?
- (2) How do we show that $\lim_{k \rightarrow \infty} x_k = \tau$?

There are some subtleties here, particularly with part (2). What we would like to say is that at each step the error is getting smaller. That is

$$|\tau - x_k| < |\tau - x_{k-1}| \quad \text{for } k = 1, 2, 3, \dots$$

But we can't. Usually all we can say is that the **bounds** on the error is getting smaller. That is: **let ε_k be a bound on the error at step k**

$$|\tau - x_k| < \varepsilon_k,$$

then $\varepsilon_{k+1} < \mu \varepsilon_k$ for some number $\mu \in (0, 1)$. It is easiest to explain this in terms of an example, so we'll study the simplest method: **Bisection**.

The most elementary algorithm is the “*Bisection Method*” (also known as “Interval Bisection”). Suppose that we know that f changes sign on the interval $[a, b] = [x_0, x_1]$ and, thus, $f(x) = 0$ has a solution, τ , in $[a, b]$. Proceed as follows

1. Set x_2 to be the midpoint of the interval $[x_0, x_1]$.
2. Choose one of the sub-intervals $[x_0, x_2]$ and $[x_2, x_1]$ where f change sign;
3. Repeat Steps 1–2 on that sub-interval, until f sufficiently small at the end points of the interval.

This may be expressed more precisely using some *pseudocode*.

Method (Bisection)

Set ϵ to be the stopping criterion.

If $|f(a)| \leq \epsilon$, return a . Exit.

If $|f(b)| \leq \epsilon$, return b . Exit.

Set $x_0 = a$ and $x_1 = b$.

Set $x_L = x_0$ and $x_R = x_1$.

Set $k = 1$

while($|f(x_k)| > \epsilon$)

$x_{k+1} = (x_L + x_R)/2$;

 if $(f(x_L)f(x_{k+1}) < 0)$

$x_R = x_{k+1}$;

 else

$x_L = x_{k+1}$

 end if;

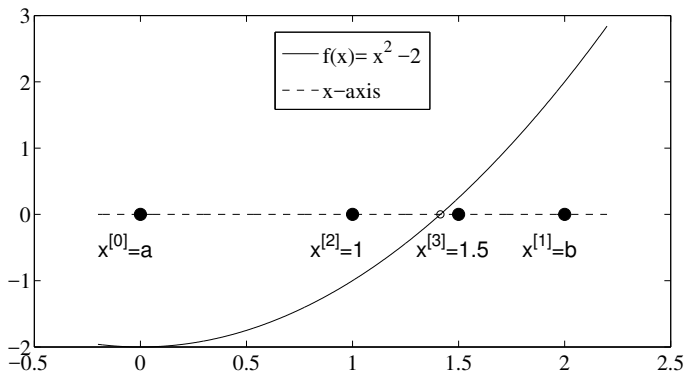
$k = k + 1$

end while;

Example

Find an estimate for $\sqrt{2}$ that is correct to 6 decimal places.

Solution: Use bisection to solve $f(x) := x^2 - 2 = 0$ on the interval $[0, 2]$.



Example

Find an estimate for $\sqrt{2}$ that is correct to 6 decimal places.

Solution: Use bisection to solve $f(x) := x^2 - 2 = 0$ on the interval $[0, 2]$.

| k | x_k | $ x_k - \tau $ | $ x_k - x_{k-1} $ |
|----------|----------|----------------|-------------------|
| 0 | 0.000000 | 1.41 | |
| 1 | 2.000000 | 5.86e-01 | |
| 2 | 1.000000 | 4.14e-01 | 1.00 |
| 3 | 1.500000 | 8.58e-02 | 5.00e-01 |
| 4 | 1.250000 | 1.64e-01 | 2.50e-01 |
| 5 | 1.375000 | 3.92e-02 | 1.25e-01 |
| 6 | 1.437500 | 2.33e-02 | 6.25e-02 |
| 7 | 1.406250 | 7.96e-03 | 3.12e-02 |
| 8 | 1.421875 | 7.66e-03 | 1.56e-02 |
| 9 | 1.414062 | 1.51e-04 | 7.81e-03 |
| 10 | 1.417969 | 3.76e-03 | 3.91e-03 |
| \vdots | \vdots | \vdots | \vdots |
| 22 | 1.414214 | 5.72e-07 | 9.54e-07 |

The main advantages of the Bisection method are

- It will always work.
- After k steps we know that

Theorem

$$|\tau - x_k| \leq \left(\frac{1}{2}\right)^{k-1} |b - a|, \quad \text{for } k = 2, 3, 4, \dots$$

A disadvantage of bisection is that it is not as efficient as some other methods we'll investigate later.

The bisection method is not very efficient. Our next goals will be to derive better methods, particularly the **Secant Method** and **Newton's method**. We also have to come up with some way of expressing what we mean by “better”; and we'll have to use Taylor's theorem in our analyses.