

## Numerical Methods 2. Solution of nonlinear equations

### Solution by interval halving (bisection method)

Let  $f : [a, b] \rightarrow \mathbf{R}$  be continuous, assume that  $f(a) < 0$ ,  $f(b) > 0$ . Let us look for a solution of the equation

$$f(x) = 0$$

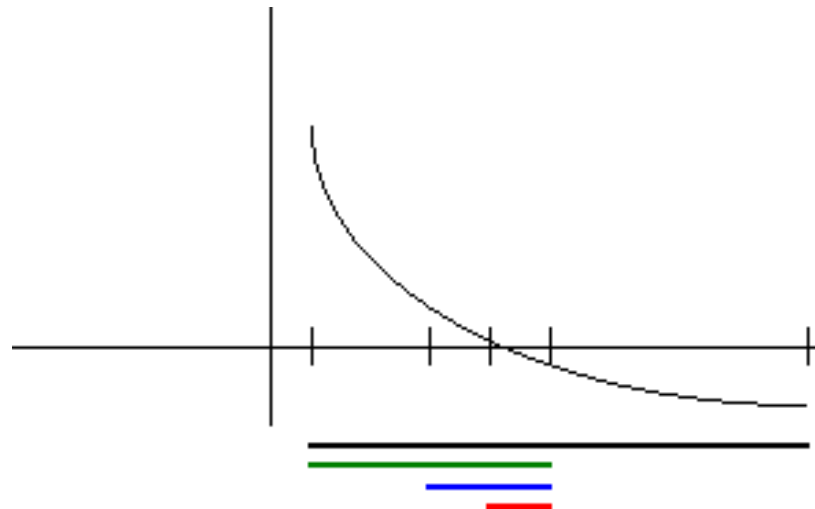
in the interval  $[a, b]$ .

**Bolzano's theorem:** If  $f$  is continuous on the finite, closed interval  $[a, b]$ , and the signs of  $f(a)$  and  $f(b)$  are different, e.g.  $f(a) < 0$ ,  $f(b) > 0$ , then  $f$  has (at least one) zero in this interval.

Let us systematically halve the interval  $[a, b]$  by taking the subinterval which has the property that the values of  $f$  at the endpoints of the subinterval have different signs. Denote by  $x_n$  the centre of the subinterval obtained in the  $n$ th step. Then this sequence converges to a zero (one of the zeroes) of the above equation. The speed of convergence is of a geometric sequence with quotient  $1/2$ .

## Solution by interval halving (bisection method)

The algorithm is illustrated by the following figure:



*Error estimation:* Defining  $x_0 := \frac{a+b}{2}$ , it is obvious that:

$$|x_n - x^*| \leq \frac{b-a}{2^n}$$

## The method based on Banach's fixed point theorem

Let  $X$  be a Banach space, let  $f : X \rightarrow X$  be a mapping, and look for a vector  $x$  such that

$$x = f(x)$$

(a **fixed point** of the mapping  $f$ ).

**Banach's fixed point theorem:** Let  $X$  be a Banach space, let  $f : X \rightarrow X$  be a **contraction** in  $X$ , i.e. assume that there exists a number  $0 \leq q < 1$  such that  $\|f(x) - f(y)\| \leq q \cdot \|x - y\|$  is valid for every  $x, y \in X$ . Then  $f$  has an unique fixed point  $x \in X$ , and this is the limit of the following, recursively defined *iteration sequence*:

$$x_0 \in X, \quad x_{n+1} := f(x_n) \quad (n = 0, 1, 2, \dots)$$

In particular, if  $f : \mathbf{R} \rightarrow \mathbf{R}$  a function for which  $\max |f'(x)| < 1$  is satisfied, then  $f$  is a contraction, since Lagrange's mean value theorem implies that

$$|f(x) - f(y)| = |f'(\xi)| \cdot |x - y| \leq (\max |f'|) \cdot |x - y|.$$

**Proof of the fixed point theorem:** The distance of two consecutive terms:

$$\begin{aligned}\|x_{n+1} - x_n\| &= \|f(x_n) - f(x_{n-1})\| \leq q \|x_n - x_{n-1}\| = q \|f(x_{n-1}) - f(x_{n-2})\| \leq q^2 \|x_{n-1} - x_{n-2}\| \leq \\ &\dots \leq q^n \|x_1 - x_0\|\end{aligned}$$

Utilizing this estimation, we show that  $(x_n)$  is a Cauchy sequence in  $X$ :

$$\begin{aligned}\|x_{n+k} - x_n\| &= \|x_{n+k} - x_{n+k-1} + x_{n+k-1} - x_{n+k-2} + \dots + x_{n+1} - x_n\| \leq \\ &\leq \|x_{n+k} - x_{n+k-1}\| + \|x_{n+k-1} - x_{n+k-2}\| + \dots + \|x_{n+1} - x_n\| \leq \\ &\leq (q^{n+k-1} + q^{n+k-2} + \dots + q^n) \cdot \|x_1 - x_0\| \leq \\ &\leq (q^n + q^{n+1} + q^{n+2} + \dots) \cdot \|x_1 - x_0\| = \frac{q^n}{1-q} \cdot \|x_1 - x_0\| \rightarrow 0\end{aligned}$$

(when  $n \rightarrow +\infty$ .) Therefore the sequence is convergent,  $x_n \rightarrow x \in X$ . We show that  $x$  is a fixed point of  $f$ . By definition:  $x_{n+1} = f(x_n)$  The left hand side obviously tends to  $x$ . The right-hand side tends to  $f(x)$ , since  $f$  is continuous. This implies that  $x = f(x)$ .

Finally, prove the uniqueness of the fixed point. If  $x, y$  were two *different* fixed points, then

$$0 < \|x - y\| = \|f(x) - f(y)\| \leq q \cdot \|x - y\| < \|x - y\|$$

would be valid, which is impossible.

## Iteration based on Banach's fixed point theorem, examples:

1. Solve the equation  $x = \frac{1}{2} \cos x$ .

The function defined by  $f(x) := \frac{1}{2} \cos x$  is a contraction in  $\mathbf{R}$ , since  $|f'(x)| = \frac{1}{2} |\sin x| \leq \frac{1}{2}$ .

Thus, an unique fixed point exists, and e.g. the following sequence converges to this:

$x_0 := 0$ ,  $x_{n+1} := \frac{1}{2} \cos x_n$ . The first terms of the sequence are as follows (with four decimal digits): 0.0000, 0.5000, 0.4387, 0.4526, 0.4496, 0.4502, 0.4501, 0.4501, 0.4501, ...

2. Let  $B \in \mathbf{M}_{N \times N}$ ,  $g \in \mathbf{R}^N$  be given, and solve the linear system of equations  $x = Bx + g$ .

If  $\|B\| < 1$  (with respect to an arbitrary matrix norm induced by a vector norm), then the mapping  $f(x) := Bx + g$  is a contraction, since

$$\|f(x) - f(y)\| = \|Bx + g - By - g\| \leq \|B\| \cdot \|x - y\|$$

In this case, there exists a unique fixed point, and the vector sequence defined by  $x_0 := \mathbf{0}$ ,  $x_{n+1} := Bx_n + g$  converges to the fixed point.

## Newton's method for univariate functions

Let  $f : (a,b) \rightarrow \mathbf{R}$  be a given function. Solve the equation

$$f(x) = 0$$

in the interval  $(a,b)$ .

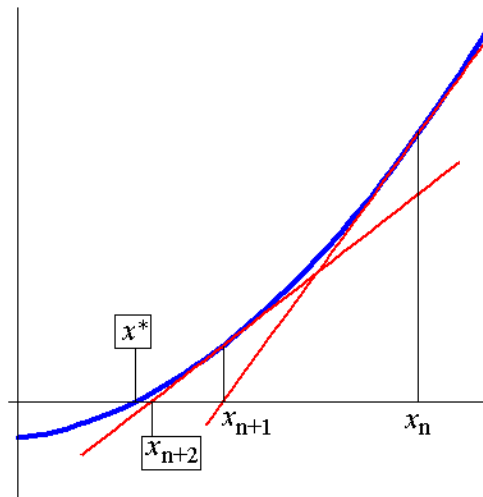
**Newton's method:** If  $x_n$  is an approximate solution, then define an improved approximation to be the zero of the tangent line at  $x_n$ . The equation of the tangent line is:

$y = f(x_n) + f'(x_n) \cdot (x - x_n)$ , whence:

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} \quad (n = 0, 1, 2, \dots) \quad (x_0 \in (a, b): \text{initial approximation})$$

## Newton's method for univariate functions

Illustration of the method ( $x^*$  denotes the exact solution):



If  $f$  is twice continuously differentiable,  $f$  has a zero in  $(a,b)$ , and  $f'(x^*) \neq 0$  is valid, then Newton's iteration **quadratically converges** to  $x^*$  for arbitrary initial approximation  $x_0$  which is sufficiently close to the exact solution  $x^*$ , i.e. for a proper positive constant  $C > 0$ :

$$|x_{n+1} - x^*| \leq C \cdot |x_n - x^*|^2$$

**Proof:** We utilize *Lagrange's mean value theorem* twice:

$$\begin{aligned} x_{n+1} - x^* &= x_n - x^* - \frac{f(x_n) - f(x^*)}{f'(x_n)} = x_n - x^* - \frac{f'(t)(x_n - x^*)}{f'(x_n)} = \\ &= \frac{f'(x_n) - f'(t)}{f'(x_n)} \cdot (x_n - x^*) = \frac{f''(s)}{f'(x_n)} \cdot (x_n - t) \cdot (x_n - x^*) \end{aligned}$$

Since  $f'(x^*) \neq 0$ , the derivative function differs from zero in a closed neighbourhood of  $x^*$ . In this neighbourhood:

$$|x_{n+1} - x^*| \leq \frac{\max |f''|}{\min |f'|} \cdot |x_n - t| \cdot |x_n - x^*| \leq \frac{\max |f''|}{\min |f'|} \cdot |x_n - x^*|^2 = C \cdot |x_n - x^*|^2$$

### Newton's method, example:

Let  $A$  be a fixed positive number, and define:  $f(x) := x^2 - A$ . ( $\Rightarrow f'(x) = 2x$ )

Now the unique positive solution of the equation

$$f(x) = 0$$

is  $x = \sqrt{A}$ .

Starting from an arbitrary initial approximation  $x_0 > 0$  (e.g.  $x_0 := A$ ), we arrive at the following recursion:

$$x_{n+1} := x_n - \frac{x_n^2 - A}{2x_n} = \frac{1}{2} \cdot \left( x_n + \frac{A}{x_n} \right)$$

The sequence converges to  $\sqrt{A}$  *extremely rapidly*, requiring only additions and divisions.

*Remark:* Newton's method can be applied to the computation of any root in a similar way.

## Some variants of Newton' method

The main difficulty: the computation of the derivatives.

**The secant method:** Here  $f'(x_n) \approx \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}$ , which results in the recursion:

$$x_{n+1} := x_n - \frac{(x_n - x_{n-1}) \cdot f(x_n)}{f(x_n) - f(x_{n-1})}$$

If  $f$  is twice continuously differentiable,  $f$  has a root  $x^*$  in  $(a,b)$ , and  $f'(x^*) \neq 0$ , then the secant method defines an iteration which converges to  $x^*$  provided that the initial approximations  $x_0, x_1$  are sufficiently close to the exact solution.

The speed of convergence is **at least** that of a geometrical sequence, i.e.

$$|x_n - x^*| \leq C \cdot q^n$$

for some  $C > 0$ ,  $0 < q < 1$

*Remark:* In fact, the speed of convergence is faster (*superlinear convergence*).

## Some variants of Newton' method

### Steffensen's method:

Assume that the function  $f : \mathbf{R} \rightarrow \mathbf{R}$  is twice continuously differentiable and has a unique root  $x^*$ . Assume also that  $f'(x^*) \neq 0$ . Then for any initial approximation  $x_0$  which is sufficiently close to  $x^*$ , the following recursive sequence is quadratically converges to  $x^*$ :

$$x_{n+1} := x_n - \frac{f(x_n)^2}{f(x_n + f(x_n)) - f(x_n)} \quad (n = 0, 1, 2, \dots)$$

*Remark:* Both the secant method and Steffensen's method require computing the values of the function  $f$  but not of the derivatives.

**Proof of the convergence of Steffensen's method:** Utilizing Lagrange's mean value theorem:

$$f(x_n + f(x_n)) - f(x_n) = f'(t) \cdot (x_n + f(x_n) - x_n) = f'(t) \cdot f(x_n)$$

therefore

$$\begin{aligned} x_{n+1} - x^* &= x_n - x^* - \frac{f(x_n)}{f'(t)} = x_n - x^* - \frac{f(x_n) - f(x^*)}{f'(t)} = x_n - x^* - \frac{f'(s)}{f'(t)}(x_n - x^*) = \\ &= \frac{f'(t) - f'(s)}{f'(t)} \cdot (x_n - x^*) = \frac{f''(w)}{f'(t)} \cdot (t - s) \cdot (x_n - x^*) \end{aligned}$$

Since  $f'(x^*) \neq 0$ , therefore the derivative function differs from a closed neighbourhood of  $x^*$ , and here:

$$|x_{n+1} - x^*| \leq \frac{\max |f''|}{\min |f'|} \cdot |t - s| \cdot |x_n - x^*| \leq \frac{\max |f''|}{\min |f'|} \cdot |x_n - x^*|^2 = C \cdot |x_n - x^*|^2$$

## Differentiation of functions mapping between Banach spaces

Let  $X, Y$  be Banach spaces. The mapping  $F : X \rightarrow Y$  is said to be **differentiable** at the point  $x \in X$  and its **derivative** is the **bounded linear operator**  $A : X \rightarrow Y$ , if for any vector  $h$  chosen from a proper neighbourhood of  $\mathbf{0}$ , the following equality is valid:

$$F(x + h) = F(x) + Ah + o(h)$$

where  $o(h)$  is an expression such that  $\frac{o(h)}{\|h\|} \rightarrow \mathbf{0} \quad (h \rightarrow \mathbf{0})$ .

Notations:  $F'(x)$  or  $DF(x)$ .

*Example:*  $F : \mathbf{R}^N \rightarrow \mathbf{R}$ ,  $F(x) := \langle Ax, x \rangle$  (where  $A \in \mathbf{M}_{N \times N}$  is a self-adjoint matrix), then:

$$F(x + h) = \langle A(x + h), x + h \rangle = \langle Ax, x \rangle + 2\langle Ax, h \rangle + \langle Ah, h \rangle = F(x) + \langle 2Ax, h \rangle + O(h^2).$$

Thus,  $F'(x) = 2Ax \in \mathbf{R}^N$ .

## Generalized Newton method

Newton's method for the equation  $F(x) = 0$ :

$$x_{n+1} = x_n - (DF(x_n))^{-1} F(x_n) \quad (n = 0, 1, 2, \dots)$$

This means that:  $x_{n+1} = x_n - w_n \quad (n = 0, 1, 2, \dots)$

where the correction term  $w_n$  is the solution of the following **linear** equation:

$$DF(x_n)w_n = F(x_n)$$

If  $F$  is twice continuously differentiable,  $F$  has a root in  $X$ , and  $DF(x^*)$  is *regular* (i.e. invertible with a bounded inverse), then Newton's method quadratically converges to the exact solution  $x^*$  provided that the initial approximation  $x_0$  is sufficiently close to  $x^*$ . That is, the following estimation is valid (with a proper constant  $C > 0$ ):

$$\|x_{n+1} - x^*\| \leq C \cdot \|x_n - x^*\|^2$$

*Remark:* Newton's method converts a **nonlinear** problem to a **sequence of linear** ones.

## Generalized Newton method, an example

**Inversion of a matrix.** Let  $A \in \mathbf{M}_{N \times N}$  be a regular matrix. For an arbitrary regular matrix  $X \in \mathbf{M}_{N \times N}$ , define the following operator:

$$F(X) := X^{-1} - A$$

Then  $F : \mathbf{M}_{N \times N} \rightarrow \mathbf{M}_{N \times N}$ , and the unique solution of the equation  $F(X) = 0$  is:  $X = A^{-1}$ .

Let us apply Newton's method to the matrix equation. First, calculate the derivative of  $F$ :

$$F(X + H) = (X + H)^{-1} - A = (X(I + X^{-1}H))^{-1} - A = (I + X^{-1}H)^{-1}X^{-1} - A$$

If the norm of the matrix  $H$  is sufficiently small, then  $\|X^{-1}H\| \leq \|X^{-1}\| \cdot \|H\| < 1$ .

Utilizing the expression  $(I - B)^{-1} = I + B + B^2 + B^3 + B^4 + \dots$  (which is valid, if  $\|B\| < 1$ , and implies that  $(I - B)^{-1} = I - B + O(\|B\|^2)$ ):

$$\begin{aligned} F(X + H) &= (I + X^{-1}H)^{-1}X^{-1} - A = (I - X^{-1}H + o(H))X^{-1} - A = \\ &= X^{-1} - X^{-1}HX^{-1} + o(H) - A = F(X) - X^{-1}HX^{-1} + o(H) \end{aligned}$$

whence

$$DF(X)H = -X^{-1}HX^{-1} \quad \Rightarrow \quad DF(X)^{-1}W = -XWX$$

## Generalized Newton method, an example

Thus, the algorithm of Newton's method is as follows:

$$\begin{aligned} X_{n+1} &:= X_n - (DF(x_n))^{-1}(X_n^{-1} - A) = X_n + X_n(X_n^{-1} - A)X_n = \\ &= X_n(2I - AX_n) \end{aligned}$$

For the error of the approximation:  $\|A^{-1} - X_n\| = \|A^{-1}(I - AX_n)\| \leq \|A^{-1}\| \cdot \|I - AX_n\|$ .

Observe that  $\|I - AX_n\|$  converges to 0 *very rapidly* (provided that the initial approximation was good enough), since:

$$I - AX_{n+1} = I - AX_n(2I - AX_n) = I - 2AX_n + AX_nAX_n = (I - AX_n)^2,$$

whence

$$\|I - AX_{n+1}\| \leq \|I - AX_n\|^2$$