

## 2. Nonlinear systems

Given a positive integer  $m$  and a continuous function  $f: \mathbb{R}^m \rightarrow \mathbb{R}^m$ , find a zero of  $f$ , i.e.,  $x \in \text{dom} f$  such that  $f(x) = 0$ .

This problem, ubiquitous in science and engineering, has been studied by some of the most brilliant minds and is still the subject of intense research nowadays. The main reason is that, in general, there is no algorithm that computes a zero of  $f$  within finitely many arithmetic operations. Thus, iterative methods are necessary.

Example. Example of a function  $f$  that has no zero:  $f(x) := x^2 + 1$ . Examples of a function  $f$  whose zeros cannot be computed within finitely many arithmetic operations: most polynomial functions of degree at least five, e.g.,  $f(x) := x^5 - 4x - 2$ , and many other functions, e.g.,  $f(x) := \cos x - x$  and  $f(x) := e^x - x - 2$ .

The general approach in numerical analysis to tackle such problems is establishing conditions that ensure the existence of a solution and designing an iterative method, i.e., a method that, given an initial guess, generates a sequence that hopefully converges to a solution.

- The two are related: the existence theorems that we are going to see in this lecture admit a constructive proof that shows that an iterative method converges to a solution.
- The speed of convergence is a criterion to compare iterative methods: the faster the convergence the better.
- Another criterion is the computational cost per iteration.

Contents:

1. bisection method and intermediate value theorem;
2. fixed-point method and Banach fixed-point theorem;
3. Newton's method;
4. secant method;
5. stopping criteria.

## 1. Bisection method and intermediate value theorem

In this section, we focus on the case where  $n=1$ . Let  $a, b \in \mathbb{R}$  such that  $a < b$ . Let  $f: \mathbb{R} \rightarrow \mathbb{R}$  be defined and continuous on  $[a, b]$ .

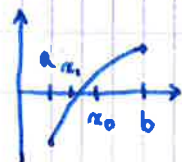
Theorem. If  $f(a)f(b) < 0$ , then there exists  $c \in (a, b)$  such that  $f(c) = 0$ .

The proof simply establishes the convergence of the bisection method. Without loss of generality, assume that  $f(a) < 0$  and  $f(b) > 0$ .

```
a0 ← a; b0 ← b; x0 ←  $\frac{a+b}{2}$ ; i ← 0;
while f(xi) ≠ 0 do
  if f(xi) < 0 then
    ai+1 ← xi; bi+1 ← bi;
  else
    ai+1 ← ai; bi+1 ← xi;
  end
  i ← i+1; xi ←  $\frac{a_i+b_i}{2}$ ;
end
```

Bisection method

Either the method finds a zero of  $f$  after a finite number of iterations or it generates infinite sequences  $(a_i)_{i \in \mathbb{N}}$ ,  $(b_i)_{i \in \mathbb{N}}$ , and  $(x_i)_{i \in \mathbb{N}}$  such that, for all  $i \in \mathbb{N}$ ,  $a_i < x_i < b_i$  and  $f(a_i) < 0$  and  $f(b_i) > 0$ . Since  $(a_i)_{i \in \mathbb{N}}$  is monotonically nondecreasing and bounded from above,  $(a_i)_{i \in \mathbb{N}}$  converges to  $\sup_{i \in \mathbb{N}} a_i \in [a, b]$ . Similarly, since  $(b_i)_{i \in \mathbb{N}}$  is monotonically nonincreasing and bounded from below,  $(b_i)_{i \in \mathbb{N}}$  converges to  $\inf_{i \in \mathbb{N}} b_i \in [a, b]$ . Moreover, since, for all  $i \in \mathbb{N}$ ,  $b_{i+1} - a_{i+1} = \frac{1}{2}(b_i - a_i)$  and thus  $b_i - a_i = 2^{-i}(b - a)$ , the sequences  $(a_i)_{i \in \mathbb{N}}$  and  $(b_i)_{i \in \mathbb{N}}$  have the same limit. Letting  $i$  tend to infinity in  $f(a_i) < 0$  and  $f(b_i) > 0$  shows that the limit is a zero of  $f$  since  $f$  is continuous.



Unfortunately, no easy extension to the case where  $n > 1$ .

## 2. Fixed-point method and Banach fixed-point theorem

The following is a direct consequence of the preceding theorem.

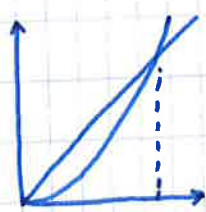
Theorem (Brouwer fixed-point theorem).

Let  $g: \mathbb{R} \rightarrow \mathbb{R}$  be defined and continuous on  $[a, b]$ . If  $g([a, b]) \subseteq [a, b]$ , then there exists  $c \in [a, b]$  such that  $g(c) = c$ .

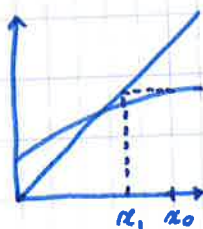
Proof. Apply the preceding theorem to the function  $f$  defined by  $f(x) := g(x) - x$  for all  $x \in [a, b]$ . ■

Definition. A fixed point of a function  $g: \mathbb{R}^n \rightarrow \mathbb{R}^n$  is a point  $x \in \text{dom } g$  such that  $g(x) = x$ .

The proof of the preceding theorem shows that the problem of finding a zero of a function  $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$  can always be reformulated as the problem of finding a fixed point of a function  $g: \mathbb{R}^n \rightarrow \mathbb{R}^n$ . Furthermore, to find a fixed point of a function  $g: \mathbb{R}^n \rightarrow \mathbb{R}^n$ , it is natural to consider the iteration  $x_{i+1} := g(x_i)$  for all  $i \in \mathbb{N}$ . Indeed, if  $g$  is continuous, then the generated sequence can converge only to fixed points of  $g$ .



no convergence  
from this zone



For convergence, the slope must be smaller than 1.

We are going to establish a sufficient condition for convergence based on the following

Definition. Let  $g: \mathbb{R}^n \rightarrow \mathbb{R}^n$  be defined on a nonempty subset  $S$  of  $\mathbb{R}^n$ . The function  $g$  is said to be Lipschitz continuous on  $S$  if

$$\text{Lip}_S(g) := \sup_{\substack{x, y \in S \\ x \neq y}} \frac{\|g(x) - g(y)\|}{\|x - y\|} < \infty.$$

The function  $g$  is called a contraction on  $S$  if  $\text{Lip}_S(g) < 1$ .  
↳ distance between  $g(x)$  and  $g(y)$  < distance between  $x$  and  $y$

Proposition. Let  $U$  be an open convex subset of  $\mathbb{R}^n$  and let  $g: U \rightarrow \mathbb{R}^n$  be differentiable. Then,  $g$  is Lipschitz continuous if and only if  $g'$  is bounded.

The proof shows that  

$$\text{Lip}(g) = \sup_U \|g'\|$$

Proof. Assume that  $g'$  is bounded. Let  $x$  and  $y$  be two distinct points in  $U$ . By the mean value theorem,

$$\|g(x) - g(y)\| \leq \|x - y\| \sup_{t \in (0,1)} \|g'(x + t(y-x))\|.$$

Thus,  $\text{Lip}(g) \leq \sup_{x \in U} \|g'(x)\|$ .

Conversely, assume that  $g$  is Lipschitz continuous. Let  $x \in U$  and  $h \in \mathbb{R}^n \setminus \{0\}$ . For all  $t \in (0, \infty)$  sufficiently small,  $x + th \in U$  and

$$\frac{\|g'(x)h\|}{\|h\|} = \frac{\|g'(x)th\|}{\|th\|} \leq \underbrace{\frac{\|g(x+th) - g(x) - g'(x)th\|}{\|th\|}}_{\rightarrow 0 \text{ as } t \rightarrow 0} + \underbrace{\frac{\|g(x+th) - g(x)\|}{\|th\|}}_{\leq \text{Lip}(g)}.$$

Thus,  $\|g'(x)\| \leq \text{Lip}(g)$ . ■

Theorem (Banach fixed-point theorem).

Let  $g: \mathbb{R}^n \rightarrow \mathbb{R}^n$  be defined on a nonempty closed subset of  $\mathbb{R}^n$ . If  $g$  is a contraction on  $C$  and  $g(C) \subseteq C$ , then  $g$  has a unique fixed point  $x_* \in C$  and, for every  $x_0 \in C$ ,  $\|g^i(x_0) - x_*\| \leq \left(\frac{\text{Lip}(g)}{c}\right)^i \|x_0 - x_*\| \forall i \in \mathbb{N}$ , which implies that  $(g^i(x_0))_{i \in \mathbb{N}}$  converges to  $x_*$ .

Proof. ① Uniqueness of the fixed point. If  $x$  and  $y$  are fixed points of  $g$ , then  $\|x - y\| = \|g(x) - g(y)\| \leq \text{Lip}(g) \|x - y\|$ , thus  $(1 - \text{Lip}(g)) \|x - y\| \leq 0$ , hence  $\|x - y\| \leq 0$ , and therefore  $x = y$ .

② Inequality. Let  $x_* \in C$  be a fixed point of  $g$ . Let  $x_0 \in C$ . For all  $i \in \mathbb{N}$ ,  $\|g^i(x_0) - x_*\| = \|g^i(x_0) - g^i(x_*)\| \leq \left(\frac{\text{Lip}(g)}{c}\right)^i \|x_0 - x_*\|$ .

③ Existence of the fixed point. Let  $L := \text{Lip}(g)$ . First, for all  $x, y \in C$ ,

$$\|x - y\| \leq \frac{\|x - g(x)\| + \|y - g(y)\|}{1 - L}$$

since

$$\|x - y\| \leq \|x - g(x)\| + \|g(x) - g(y)\| + \|g(y) - y\| \leq \|x - g(x)\| + L \|x - y\| + \|g(y) - y\|.$$

Based on this inequality, let us prove that, for every  $x_0 \in C$ ,  $(g^i(x_0))_{i \in \mathbb{N}}$  is a Cauchy sequence. For all  $i, j \in \mathbb{N}$ ,

$$\|g^i(x_0) - g^j(x_0)\| \leq \frac{\|g^i(x_0) - g^{i+1}(x_0)\| + \|g^j(x_0) - g^{j+1}(x_0)\|}{1 - L} \leq \frac{L^i + L^j}{1 - L} \|x_0 - x_*\|.$$

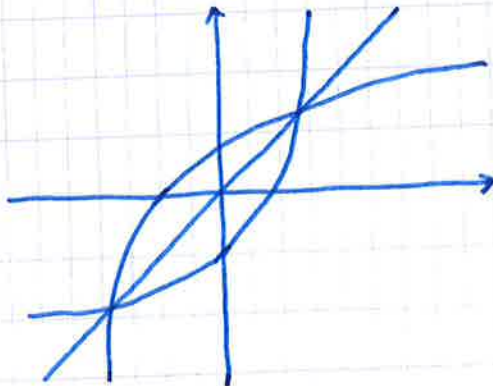
Since  $\mathbb{R}^m$  is complete,  $(g^i(\alpha_0))_{i \in \mathbb{N}}$  converges. Moreover, letting  $i$  tend to infinity in  $g^{i+1}(\alpha_0) = g(g^i(\alpha_0))$  shows that the limit is a fixed point of  $g$  since  $g$  is continuous on  $C$ . ■

There are often several ways of turning the problem of finding a zero into that of finding a fixed point, and some ways may be better than others, as illustrated next.

Example. Let  $f: \mathbb{R} \rightarrow \mathbb{R}: x \mapsto e^x - x - 2$ . Then,  $f(-2) > 0$ ,  $f(-1) < 0$ ,  $f(1) < 0$ , and  $f(2) > 0$ . Thus,  $f$  has a zero on  $(-2, -1)$  and a zero on  $(1, 2)$ .

- Define  $g: \mathbb{R} \rightarrow \mathbb{R}: x \mapsto e^x - 2$ . Then, the zeros of  $f$  are exactly the fixed points of  $g$ . For all  $x \in (-2, -1)$ ,  $-2 < g(-2) < g(x) < g(-1) < -1$ . For all  $x \in \mathbb{R}$ ,  $g'(x) = e^x$ . Thus, for all  $x \in (-2, -1)$ ,  $0 < e^{-2} = g'(-2) < g'(x) < g'(-1) = e^{-1} < 1$ . Hence,  $g$  is a contraction on  $[-2, -1]$ . However, for all  $x \in [0, \infty)$ ,  $1 \leq g'(x)$ . Therefore,  $g$  is not a contraction on any subinterval of  $[0, \infty)$ . In conclusion, with this  $g$ , the fixed-point method can be used to find the fixed point in  $(-2, -1)$  but not the fixed point in  $(1, 2)$ .

- Define  $g: (-2, \infty) \rightarrow \mathbb{R}: x \mapsto \ln(x+2)$ . For all  $x \in (1, 2)$ ,  $1 < \ln 3 = g(1) < g(x) < g(2) = 2 \ln 2 < 2$ . For all  $x \in (-2, \infty)$ ,  $g'(x) = \frac{1}{x+2}$ . Thus, for all  $x \in (1, 2)$ ,  $0 < \frac{1}{4} = g'(2) < g'(x) < g'(1) = \frac{1}{3} < 1$ . Hence,  $g$  is a contraction on  $[1, 2]$ . However, for all  $x \in (-2, -1)$ ,  $1 \leq g'(x)$ . Therefore,  $g$  is not a contraction on  $(-2, -1)$ . In conclusion, with this  $g$ , the fixed-point method can be used to find the fixed point in  $(1, 2)$  but not the fixed point in  $(-2, -1)$ .



### 3. Newton's method

Faster convergence can be guaranteed if  $f$  is differentiable. The idea is to replace  $f$  with its first-order Taylor polynomial at every iteration. Thus, given  $\alpha_0 \in \text{dom } f$ , solve  $f(\alpha_0) + f'(\alpha_0)(\alpha - \alpha_0) = 0$  instead of  $f(\alpha) = 0$ . If  $f'(\alpha_0)$  is invertible, then the linear system has a unique solution  $\alpha_1 := \alpha_0 - f'(\alpha_0)^{-1} f(\alpha_0)$ . This yields the following method.

```
i ← 0
while f(αi) ≠ 0 do
  try to find α ∈ ℝn such that f'(αi)α = -f(αi), e.g., by Gaussian el.
  if f'(αi) is not invertible then
    stop
  else
    αi+1 ← αi + α
    i ← i+1
  end
end
```

### Newton's method

Newton's method is a fixed-point method for  $g: \mathbb{R}^n \rightarrow \mathbb{R}^n: \alpha \mapsto \alpha - f'(\alpha)^{-1} f(\alpha)$ .

Theorem (Ortega & Rheinboldt 2000, § 10.2.2).

Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$  be defined and differentiable on a nonempty open subset  $U$  of  $\mathbb{R}^n$ . Let  $\alpha_* \in U$  be a zero of  $f$ . If  $f'$  is continuous at  $\alpha_*$  and  $f'(\alpha_*)$  is invertible, then there exists  $\delta \in (0, \infty)$  such that  $B(\alpha_*, \delta) \subseteq U$ ,  $f'(\alpha)$  is invertible for all  $\alpha \in B(\alpha_*, \delta)$ , and  $\lim_{\alpha \rightarrow \alpha_*} \frac{\|g(\alpha) - \alpha_*\|}{\|\alpha - \alpha_*\|} = 0$ , which implies that for all  $\alpha \in B(\alpha_*, \delta)$  sufficiently close to  $\alpha_*$ , the sequence  $(g^i(\alpha))_{i \in \mathbb{N}}$  converges to  $\alpha_*$ . If, moreover, there exist  $L \in [0, \infty)$  and  $p \in (0, 1]$  such that, for all  $\alpha \in B(\alpha_*, \delta)$ ,  $\|f'(\alpha) - f'(\alpha_*)\| \leq L \|\alpha - \alpha_*\|^p$ , then  $\|g(\alpha) - \alpha_*\| \leq 4L \|f'(\alpha_*)^{-1}\| \|\alpha - \alpha_*\|^{p+1}$  for all  $\alpha \in B(\alpha_*, \delta)$ .

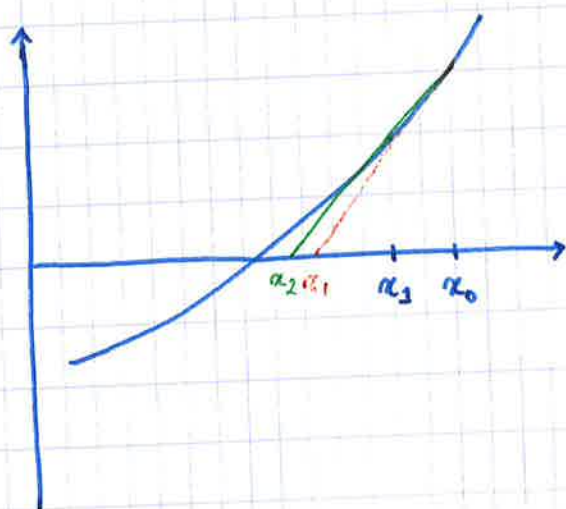
#### 4. Secant method

In this section, we focus on the case where  $n=1$ . To avoid computing the derivative of  $f$ , we can replace  $f$  with its secant instead of its tangent.

Thus, given  $\alpha_0, \alpha_1 \in \text{dom } f$ ,  $\alpha_0 \neq \alpha_1$ , we solve  $f(\alpha_0) + \frac{f(\alpha_1) - f(\alpha_0)}{\alpha_1 - \alpha_0} (\alpha - \alpha_0) = 0$  instead of  $f(\alpha) = 0$ . If  $f(\alpha_0) \neq f(\alpha_1)$ , then the unique solution is  $\alpha_2 := \frac{\alpha_0 f(\alpha_1) - \alpha_1 f(\alpha_0)}{f(\alpha_1) - f(\alpha_0)}$ . This yields the following method.

```
i ← 0
while 0 ≠ f(αi) ≠ f(αi+1) ≠ 0 do
  αi+2 ←  $\frac{\alpha_i f(\alpha_{i+1}) - \alpha_{i+1} f(\alpha_i)}{f(\alpha_{i+1}) - f(\alpha_i)}$ 
  i ← i+1
end
```

Secant method



Newton  
secant

#### 5. Stopping criteria

Stop when  $\|f(\alpha_i)\| \leq \epsilon$  or  $\|\alpha_{i+1} - \alpha_i\| \leq \epsilon$  for some  $\epsilon \in (0, \infty)$ . The inequalities obtained in the convergence analyses sometimes enable to find  $i \in \mathbb{N}$  such that the second stopping criterion is satisfied. For illustration, let us find  $i \in \mathbb{N}$  such that  $\|\alpha_i - \alpha_*\| \leq \epsilon$  for the bisection method, Newton's method, and the fixed-point method.

• Bisection method:  $\left\lceil \log_2 \left( \frac{b-a}{\epsilon} \right) - 1 \right\rceil$

For all  $i \in \mathbb{N}$ ,  $|x_i - x_*| \leq \frac{1}{2}(b_i - a_i) = 2^{-i-1}(b-a)$  thus  $|x_i - x_*| \leq \epsilon$  if  $2^{-i-1}(b-a) \leq \epsilon$ , i.e.,  $i \geq \log_2 \left( \frac{b-a}{\epsilon} \right) - 1$ .

• Fixed-point method:  $\left\lceil \frac{\ln \left( \frac{\epsilon(1-Lip(g))}{\|x_0 - x_*\|} \right)}{\ln Lip(g)} \right\rceil$

For all  $i \in \mathbb{N}$ ,  $\|x_i - x_*\| \leq \left( Lip(g) \right)^i \|x_0 - x_*\|$  thus  $\|x_i - x_*\| \leq \epsilon$  if  $\left( Lip(g) \right)^i \|x_0 - x_*\| \leq \epsilon$ , i.e.,  $i \geq \frac{\ln \left( \frac{\epsilon}{\|x_0 - x_*\|} \right)}{\ln Lip(g)}$ .

Moreover,  $\|x_0 - x_*\| \leq \|x_0 - x_1\| + \|x_1 - x_*\|$  thus  $\|x_0 - x_*\| \leq \frac{1}{1-Lip(g)} \|x_0 - x_1\|$ .

Hence,  $\|x_i - x_*\| \leq \epsilon$  if  $\frac{\left( Lip(g) \right)^i}{1-Lip(g)} \|x_0 - x_1\| \leq \epsilon$ , i.e.,

$$i \geq \frac{\ln \left( \frac{\epsilon(1-Lip(g))}{\|x_0 - x_1\|} \right)}{\ln Lip(g)}$$

• Newton's method in the case where  $p=1$ :  $\left\lceil \log_2 \left( \frac{\ln(c(x_*)\epsilon)}{\ln(c(x_*)\|x_0 - x_*\|)} \right) \right\rceil$

For all  $x \in B(x_*, \delta)$ ,  $\|g(x) - x_*\| \leq c(x_*) \|x - x_*\|^2$  with  $c(x_*) := 4L \|f'(x_*)^{-1}\|$ , thus  $g(x) \in B(x_*, \delta)$  if  $c(x_*) \delta \|x - x_*\| \leq \delta$ , i.e.,  $\|x - x_*\| \leq \frac{1}{c(x_*)}$ .

Let  $e := \min \left\{ \delta, \frac{1}{c(x_*)} \right\}$ . Then,  $g(B(x_*, e)) \subseteq B(x_*, e)$ .

Thus, for all  $x_0 \in B(x_*, e)$  and  $i \in \mathbb{N}$ ,  $\|x_i - x_*\| \leq \frac{1}{c(x_*)} \left( c(x_*) \|x_0 - x_*\| \right)^{2^i}$ , hence  $\|x_i - x_*\| \leq \epsilon$  if  $\frac{1}{c(x_*)} \left( c(x_*) \|x_0 - x_*\| \right)^{2^i} \leq \epsilon$ , i.e.,

$$i \geq \log_2 \left( \frac{\ln(c(x_*)\epsilon)}{\ln(c(x_*)\|x_0 - x_*\|)} \right)$$

### References:

- J.M. Ortega & W.C. Rheinboldt, Iterative Solution of Nonlinear Equations in Several Variables, SIAM, 2000;
- E. Süli & D.F. Mayers, An Introduction to Numerical Analysis, Cambridge University Press, 2012.