

# 1

## II. Nonlinear Equations

Goals: - solve nonlinear (systems of) equations numerically  
- understand that this can be hard...

Why? - Generally, the nonlinear equations that appear in practice are not solvable by analytical means  $x = \dots$

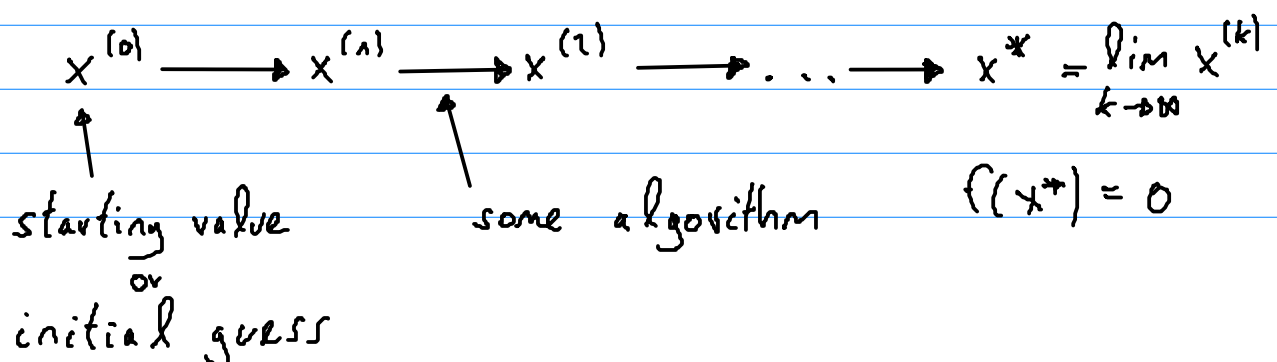
- Implicit methods for ODEs ... Very important for stiff problems (... very common in practice)
- Steady state CSTR example

MATLAB: `fsolve`

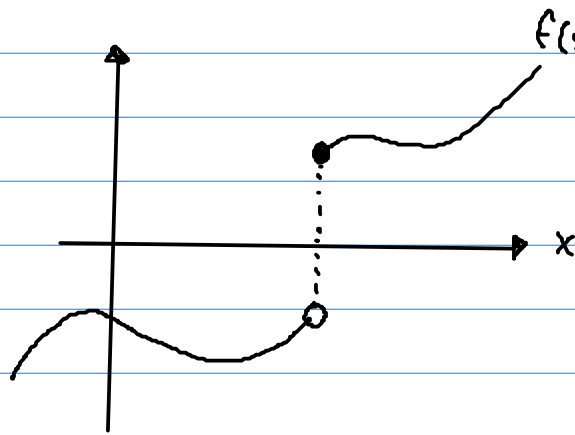
### II.1 Single nonlinear equations

Problem: Solve  $f(x) = 0$  for  $f: D=[a,b] \rightarrow \mathbb{R}$

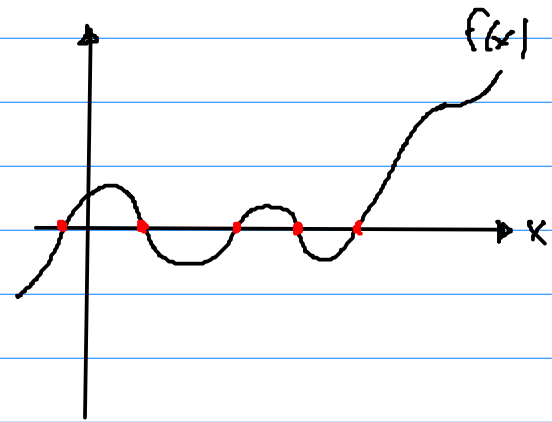
Iterative techniques



Remark: In general no existence nor uniqueness



$f$  not continuous



several solutions

Assume  $f$  continuous in the following...

## II.1.1 Fixed-point iterations and a few generalities

Our problem  $f(x) = 0$  can be written as

$$x = \phi(x)$$

for some  $\phi(x)$ . This is known as a fixed-point equation (FPE) and  $x^*$  satisfying

$$x^* = \phi(x^*)$$

as a fixed-point (FP).

Def.: A fixed-point iteration (FPI) is defined as

$$x^{(k+1)} = \phi(x^{(k)}), \quad k = 0, 1, 2, \dots$$

Def.: A fixed-point equation is said to be consistent with the original equation if

$$f(x^*) = 0 \iff x^* = \phi(x^*)$$

Ex.: (1) Solve  $f(x) = x \cdot e^x - 1 = 0$  for  $x \in [0, 1]$

(i)  $x \cdot e^x - 1 = 0$

$$x \cdot e^x = 1$$

$$x = e^{-x} = \phi_1(x) \quad \text{fixed-point eq.} \\ \text{(consistent \checkmark)}$$

(ii)  $x = \phi_2(x)$  with  $\phi_2(x) = \frac{x^2 \cdot e^x + 1}{e^x(1+x)}$

(iii)  $x = \phi_3(x)$  with  $\phi_3(x) = x + 1 - x \cdot e^x$

→ slides

... Exercise: show consistency!

Rem.: (i) Fixed-point iterations not unique

(ii) Fixed-point iterations may not converge

(iii) If they converge, they may do that with different speeds

Def.: A sequence  $x^{(k)}$  with limit  $x^*$  converges with order  $p \geq 1$ , if there exists a constant  $C > 0$  such that

$$|x^{(k+1)} - x^*| \leq C |x^{(k)} - x^*|^p$$

for all sufficiently large  $k$ .

For  $p=1$ , it must  $0 < C < 1$ .

The constant  $C$  is called the rate of convergence.

In particular, convergence with order  $\begin{cases} p=1 \\ p=2 \end{cases}$

is called  $\begin{cases} \text{linear} \\ \text{quadratic} \end{cases}$ .

It is often helpful, e.g. for code verification, to measure  $C$  and  $p$  in numerical experiments.

For this we define the error at the  $k$ -th iteration as

$$E^{(k)} = |x^{(k)} - x^*|$$

Then we can write

$$E^{(k)} = C \cdot (E^{(k-1)})^p$$

$$E^{(k+1)} = C \cdot (E^{(k)})^p$$

Taking the log on both sides gives

$$\log(E^{(k)}) = \log(C) + p \cdot \log(E^{(k-1)})$$

$$\log(E^{(k+1)}) = \log(C) + p \cdot \log(E^{(k)})$$

This can be solved for  $C$  and  $p$ :

$$p = \frac{\log(E^{(k+1)}) - \log(E^{(k)})}{\log(E^{(k)}) - \log(E^{(k-1)})}$$

$$C = \frac{E^{(k+1)}}{(E^{(k)})^p} = \frac{E^{(k)}}{(E^{(k-1)})^p}$$

Ex.: (2) Determine  $C$  &  $p$  for the fixed-point iterations from Ex. (1) ~~no~~ slides

Catch: How to compute  $E^{(k)} = |x^{(k)} - x^*|$  ?  
= ?

This brings us to the practical question:

When to stop iterating?

Stopping criteria (SC):

$$(SC1) \quad |x^{(k)} - x^{(k-1)}| \leq \text{atol} \quad (\text{absolute-})$$

$$(SC2) \quad |x^{(k)} - x^{(k-1)}| \leq \text{rtol} \cdot |x^{(k)}| \quad (\text{relative-})$$

$$(SC3) \quad |x^{(k)} - x^{(k-1)}| \leq \text{tol} \cdot (1 + |x^{(k)}|) \quad (\text{hybrid-})$$

$$(SC4) \quad |f(x^{(k)})| \leq \text{ftol} \quad (\text{function-})$$

tolerance

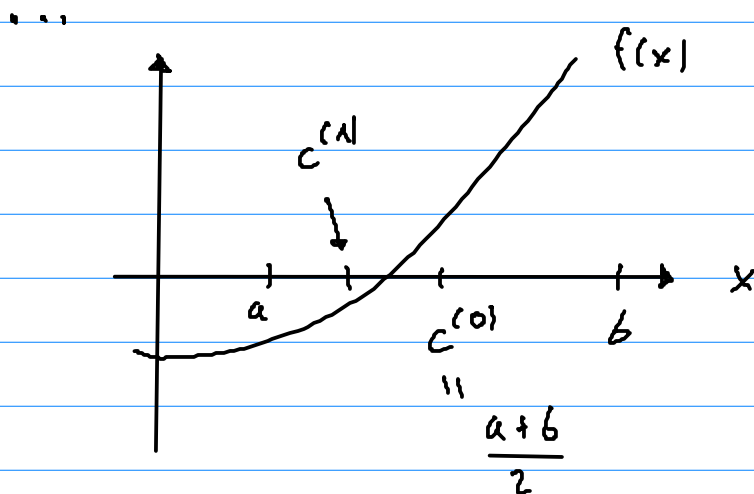
## II.1.2 Bisection method

Assume: we know  $a, b$  with

$$f(a) < 0 \quad \text{and} \quad f(b) > 0$$

$\leadsto$  root(s) bracketed

Idea: divide / bisect interval and keep subinterval fulfilling above assumption



Rem.: (i) Very easy and robust (need only  $f$ )

(ii) A priori error estimate  $\epsilon^{(k)} \leq \frac{b-a}{2^{k+1}}$

(iii) Slow convergence (linear)

(iv) Not easy to generalize to systems

## II.1.3 Newton's method

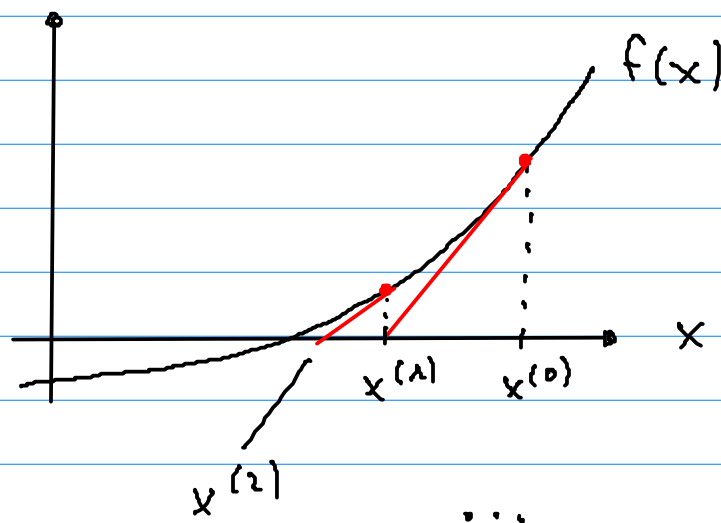
Idea: linearize  $f$

Taylor at  $x^{(k)}$ :

$$f(x) = \underbrace{f(x^{(k)}) + f'(x^{(k)}) \cdot (x - x^{(k)}) + \frac{1}{2} f''(x^{(k)}) (x - x^{(k)})^2 + \dots}_{\tilde{f}(x) \stackrel{!}{=} 0 \rightsquigarrow x^{(k+1)}}$$

$$\rightsquigarrow \tilde{f}(x^{(k+1)}) = f(x^{(k)}) + f'(x^{(k)}) \cdot (x^{(k+1)} - x^{(k)}) = 0$$

$$\rightsquigarrow x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})}$$





Rem.: (i) Needs  $f'$

(ii) Quadratic convergence ( $p=2$ ) when close enough

(iii) Can fail, e.g.  $f'(x^{(k)}) = 0$

(iv) Generalizes to systems

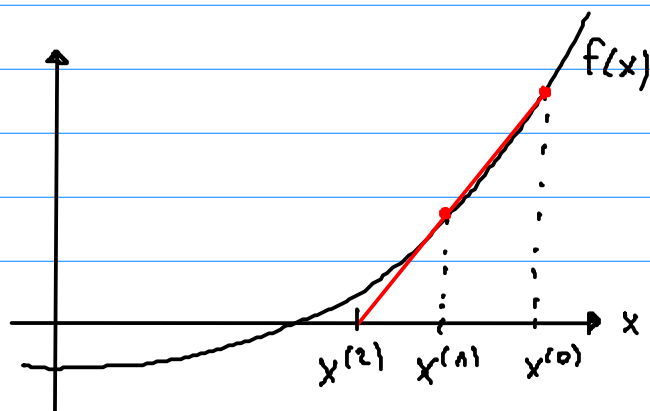
## II.1.4 Secant method

Idea: If  $f'$  is not directly available or expensive to compute, approx.  $f'$  with finite differences

$$f'(x^{(k)}) \approx \frac{f(x^{(k)}) - f(x^{(k-1)})}{x^{(k)} - x^{(k-1)}}$$

Then

$$x^{(k+1)} = x^{(k)} - f(x^{(k)}) \cdot \frac{x^{(k)} - x^{(k-1)}}{f(x^{(k)}) - f(x^{(k-1)})}$$



Rem.: (i) No need for  $f'$ !

... so-called quasi-Newton methods

(ii) Converges with  $\rho = \frac{1}{2}(1 + \sqrt{5}) \approx 1.618$

... faster than bisection, but slower than Newton

(iii) Fails when  $f(x^{(k)}) = f(x^{(k-1)})$

(iv) Generalizes to systems

## II.2 Systems of nonlinear equations

Problem: Given an  $n$ -dimensional real and continuous function  $\vec{f}: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ , solve

$$\vec{f}(\vec{x}) = 0$$

Short notation for system of nonlinear equations:

$$\left. \begin{array}{l} f_1(x_1, \dots, x_n) = 0 \\ f_2(x_1, \dots, x_n) = 0 \\ \vdots \\ f_n(x_1, \dots, x_n) = 0 \end{array} \right\} \vec{f}(\vec{x}) = 0$$

Ex.: (3)  $n=2$ ,  $D = [0, 2]^2 \subset \mathbb{R}^2$

$$f_1(x_1, x_2) = x_1^2 + x_2 - 2 = 0$$

$$f_2(x_1, x_2) = x_2 \cdot e^{x_1} - 2 = 0$$

~ contour / level curves (slides)

## II.2.1 Newton's method

Idea: linearize  $\vec{f}$

Taylor at  $\vec{x}^{(k)}$  Jacobian Matrix  $D\vec{f}(\vec{x}) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_n}{\partial x_1} & \dots & \frac{\partial f_n}{\partial x_n} \end{pmatrix}$

$$\vec{f}(\vec{x}) = \vec{f}(\vec{x}^{(k)}) + D\vec{f}(\vec{x}^{(k)}) (\vec{x} - \vec{x}^{(k)}) + \dots$$

⏟

$$\vec{f} \approx \vec{f}(\vec{x}^{(k)}) + D\vec{f}(\vec{x}^{(k)}) (\vec{x} - \vec{x}^{(k)}) \stackrel{!}{=} 0$$

$$\rightsquigarrow \vec{x}^{(k+1)} = \vec{x}^{(k)} - D\vec{f}(\vec{x}^{(k)})^{-1} \vec{f}(\vec{x}^{(k)})$$



inverse of Jacobian

Ex.:(4) For example (3)

$$D\vec{f}(\vec{x}) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} = 2x_1 & \frac{\partial f_1}{\partial x_2} = 1 \\ \frac{\partial f_2}{\partial x_1} = x_2 \cdot e^{x_1} & \frac{\partial f_2}{\partial x_2} = e^{x_1} \end{pmatrix} = \begin{pmatrix} 2x_1 & 1 \\ x_2 \cdot e^{x_1} & e^{x_1} \end{pmatrix}$$

$$D\vec{f}^{-1}(\vec{x}) = \frac{1}{(2x_1 - x_2) \cdot e^{x_1}} \begin{pmatrix} e^{x_1} & -1 \\ -x_2 \cdot e^{x_1} & 2x_1 \end{pmatrix}$$

→ slides

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

$$A^{-1} = \frac{1}{ad - cb} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$$

(5) CSTR

→ slides

Rem.: (i) Needs Jacobian

(ii) Quadratic convergence when close enough

(iii) Fails when Jacobian singular

i.e.  $D\vec{f}$  not  
invertible

...  $\sim f'(x) = 0$  ...